



# Performances statistiques d'estimateurs non-linéaires

Michael Chichignoud

## ► To cite this version:

Michael Chichignoud. Performances statistiques d'estimateurs non-linéaires. Mathématiques [math]. Université de Provence - Aix-Marseille I, 2010. Français. NNT : . tel-00540963

**HAL Id: tel-00540963**

**<https://theses.hal.science/tel-00540963>**

Submitted on 29 Nov 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

UNIVERSITÉ DE PROVENCE  
U.F.R. M.I.M.  
ÉCOLE DOCTORALE DE MATHÉMATIQUES ET INFORMATIQUE E.D. 184

## THÈSE

présentée pour obtenir le grade de  
DOCTEUR DE L'UNIVERSITÉ DE PROVENCE  
*Spécialité : Mathématiques*

par  
**Michaël CHICHIGNOUD**

sous la direction du Pr. Oleg LEPSKI

*Titre :*

**Performances statistiques d'estimateurs non-linéaires**

soutenue publiquement le 25 novembre 2010

### JURY

M. Laurent CAVALIER	Université de Provence	Examineur
Mme Béatrice LAURENT	INSA de Toulouse	Examinatrice
M. Oleg LEPSKI	Université de Provence	Directeur
M. Pascal MASSART	Université Paris-Sud	Rapporteur
Mme Dominique PICARD	Université Paris-Diderot	Examinatrice
M. Alexandre TSYBAKOV	Université Pierre et Marie Curie	Examineur
M. Aad Van der VAART	Vrije Universiteit	Rapporteur



*A mes grands parents,  
Jacqueline et Robert,  
une pensée pour eux.*



# REMERCIEMENTS

Mes premières pensées vont à mon directeur de thèse, Oleg Lepski, qui m'a fait l'immense honneur de m'accompagner pendant ces trois longues années. Il est difficile en quelques lignes d'exprimer tous mes remerciements et mon estime pour lui. J'ai souvent mis à contribution ses nombreuses qualités humaines, mais son soutien et ses conseils avisés m'ont permis de réaliser ce manuscrit. Par ses connaissances scientifiques et son extraordinaire niveau en mathématiques, qui m'impressionne toujours, il a su me lancer sur un sujet très ouvert et je le remercie très sincèrement. Dans les années futures, j'espère pouvoir répondre aux nombreuses questions ouvertes, qu'il m'a posé (et auxquelles je n'ai pas su répondre), sur des problèmes qui occuperont certainement une place importante dans les statistiques de demain.

J'exprime tous mes remerciements à Pascal Massart et Aad Van der Vaart pour avoir gentiment accepté de rapporter cette thèse. Je suis très touché qu'ils aient pris le temps de relire mes modestes travaux.

Je tiens également à remercier Laurent Cavalier, Béatrice Laurent, Dominique Picard et Alexandre Tsybakov, de l'immense honneur qu'ils me font de participer à ce jury de thèse d'une très grande qualité scientifique.

Un grand merci à tous les membres de l'équipe Probabilités/Statistiques, notamment à Florent Autin, Yuri Golubev, Grégory Maillard, Christophe Pouet et Thomas Willer pour leurs conseils avisés, leur sympathie et leur engouement pour les mathématiques.

Je voudrais remercier aussi les personnes qui ont contribué à la réalisation de cet ouvrage, notamment Christophe Pouet et Thomas Willer. Ma sincère gratitude et mon amitié vont à Joseph Salmon, futur docteur lui aussi. Un petit clin d'oeil à Adrien Saumard, qui a su, en peu de temps, me poser des questions très ouvertes sur mon sujet et me donner quelques conseils de rédaction.

Je tiens à remercier les enseignants qui m'ont le plus marqué au cours de mes études, et qui m'ont fait aimer les mathématiques : M. Fournier (Professeur de mathématiques en Première) pour son enthousiasme et son dé à *28 faces*, Dominique Barbolosi (Professeur en Licence) pour la qualité de ses cours sans support papier, et enfin Laurent Cavalier (Professeur en Master) pour ses cours dynamiques effectués avec beaucoup de rigueur.

Ces trois années passées au CMI m'ont permis de rencontrer des gens extraordinaires. Je remercie le personnel administratif pour leur disponibilité et leur sympathie, les charges

administratives étaient bien peu lourdes grâce à leur présence. Je suis heureux d'avoir fait partie de la Team des doctorants du Cmi et je remercie chacun d'entre eux pour les nombreuses discussions que nous avons eu sur divers domaines des mathématiques. Je remercie sincèrement Hamish Short (directeur de l'école doctorale) et Etienne Pardoux (ex-directeur de l'école doctorale) pour tout ce qu'ils ont fait pour moi, en particulier pour l'obtention d'une bourse doctorale. Un grand merci aux membres du bureau 114 que j'ai pu côtoyer au quotidien, notamment Clément Marteau, Sébastien Loustau et Shanti Gibert. Leur présence fût chaleureuse et indispensable durant toutes ces années.

J'exprime mes plus sincères remerciements à ma famille pour son accompagnement dans ce long parcours. Ma mère, mon frère et mes grands parents ont été d'un soutien sans faille. En particulier, un immense merci à mon père qui à force de patience et d'encouragements, m'a communiqué son penchant pour les mathématiques et la logique. Du fond du coeur Merci !

Et pour finir, tous mes sentiments vont à ma future femme Mylène, qui a accepté de me dire "oui" pour me rendre le plus heureux des hommes. Sans elle et son chaleureux soutien, je ne serais pas là aujourd'hui. Merci mon amour !

# SOMMAIRE

Le thème de cette thèse est l'estimation non-paramétrique, en particulier l'étude des performances théoriques de l'estimation de fonctions de régression, une partie importante d'un des domaines mathématiques connue sous le nom de *statistique mathématique*.

Une partie de cette thèse est rédigée en anglais, car ce sont des preprints ou des articles soumis. Les chapitres 1 et 2 sont quant à eux rédigés en français. C'est pourquoi nous avons fait précéder chaque chapitre d'un court résumé en français.

Ce manuscrit est organisé comme suit :

- Le chapitre 1 est une introduction aux modèles de régression, aux critères de performance en vigueur et aux méthodes adaptatives avec sélection de fenêtres.
- Dans le chapitre 2, le lecteur peut trouver les principaux résultats de cette thèse pour deux nouveaux types d'estimateurs. Notons qu'un lecteur non-spécialiste du domaine aura recours au chapitre 1 pour sa compréhension.  
Nous proposons aussi quelques perspectives à ce travail. Notamment, une liste d'une quinzaine de problèmes ouverts est donnée à la fin de ce chapitre.
- Les chapitres 3 et 4 traitent de l'estimation bayésienne.
- Le chapitre 5 introduit l'estimateur de Huber et quelques résultats sur ses performances.

Les résultats théoriques sont accompagnés d'expériences numériques. En particulier, on pourra comparer les estimateurs *bayésien* et de *Huber* avec les estimateurs linéaires.

Les chapitres 3 à 5 peuvent être lus indépendamment des autres (ce qui est à l'origine de quelques répétitions). Il reste néanmoins quelques liens (limités aux outils probabilistes et à des résultats techniques communs). Nous avons fait tout notre possible pour uniformiser nos notations, que nous définissons avant le chapitre 3 et qui sont redéfinies ensuite pour chaque chapitre.

La version électronique de cette thèse (à télécharger sur la page web de l'auteur : <http://www.latp.univ-mrs.fr/~chichign/doku.php?id=accueil>) comprend des liens *hyper-ref* qui permettent de se retrouver au chapitre, à la section, à la formule ou à la citation en un simple clic sur la référence. Nous incitons le lecteur à l'utiliser pour une recherche



bibliographique ou un résultat technique, bien que la version papier reste la plus agréable à lire.

# Table des matières

<b>1</b>	<b>Régression Non-Paramétrique</b>	<b>13</b>
1.1	Objet de la Thèse . . . . .	13
1.2	Modèles de Régression et Espaces Fonctionnels . . . . .	16
1.2.1	Espaces de Hölder Isotropes . . . . .	16
1.2.2	Régression Générale . . . . .	18
1.2.3	Régression Additive . . . . .	18
1.2.4	Régression Gaussienne et de Cauchy . . . . .	20
1.2.5	Régression Inhomogène de Poisson . . . . .	20
1.2.6	Régression $\alpha$ . . . . .	22
1.2.7	Régression Multiplicative Uniforme . . . . .	24
1.3	Approche Localement Paramétrique . . . . .	26
1.3.1	Estimateur Bayésien . . . . .	28
1.3.2	Estimateur de Huber . . . . .	30
1.4	Mesure de l'Erreur . . . . .	33
1.4.1	Approche Minimax Ponctuelle . . . . .	33
1.4.2	Approche Minimax Adaptative . . . . .	35
1.5	Adaptation . . . . .	36
1.5.1	Généralités . . . . .	37
1.5.2	Choix de la Fenêtre : Méthode de Lepski . . . . .	38
<b>2</b>	<b>Résultats et Perspectives</b>	<b>53</b>
2.1	Approche Bayésienne . . . . .	53
2.1.1	Recherche de la Vitesse Minimax . . . . .	54
2.1.2	Procédure Adaptative . . . . .	56

2.1.3	Grandes Déviations . . . . .	57
2.1.4	Exemples de Modèles avec des Vitesses Différentes . . . . .	58
2.2	Critère de Huber . . . . .	61
2.2.1	Adaptation . . . . .	61
2.2.2	Grandes Déviations . . . . .	63
2.2.3	Inégalités Maximales pour les processus empiriques . . . . .	65
2.3	Expériences numériques . . . . .	69
2.4	Perspectives . . . . .	79
2.4.1	Approche Bayésienne . . . . .	79
2.4.2	Critère de Huber . . . . .	81
<b>3</b>	<b>General Locally Bayesian Approach</b>	<b>85</b>
3.1	Introduction . . . . .	85
3.2	Minimax Estimation . . . . .	89
3.3	Adaptive Rule . . . . .	91
3.4	Applications . . . . .	92
3.4.1	Gaussian Regression . . . . .	93
3.4.2	Inhomogeneous Poisson Regression . . . . .	94
3.4.3	$\alpha$ Regression . . . . .	95
3.4.4	Multiplicative Uniform Regression . . . . .	96
3.5	Proofs of Main Results . . . . .	97
3.5.1	Auxiliary Results: Large Deviations . . . . .	97
3.5.2	Proof of Theorem 4 . . . . .	98
3.5.3	Proof of Theorem 5 . . . . .	98
3.5.4	Proof of Proposition 3 . . . . .	100
3.6	Appendix . . . . .	103
<b>4</b>	<b>Locally Bayesian Approach for Multiplicative Uniform Regression</b>	<b>113</b>
4.1	Introduction . . . . .	113
4.2	Minimax estimation on isotropic Hölder class . . . . .	120
4.3	Adaptive estimation on isotropic Hölder classes . . . . .	121
4.4	Simulation study . . . . .	125
4.5	Proofs of main results: upper bounds . . . . .	128

4.5.1	Auxiliary results . . . . .	128
4.5.2	Proof of Proposition 4 . . . . .	129
4.5.3	Proof of Proposition 5 . . . . .	135
4.5.4	Proof of Theorem 11 . . . . .	136
4.5.5	Proof of Theorem 14 . . . . .	137
4.6	Proofs of lower bounds . . . . .	139
4.6.1	Proof of Theorem 10 . . . . .	140
4.6.2	Proof of Theorem 13 . . . . .	140
4.6.3	Proof of Proposition 6 . . . . .	141
4.7	Appendix . . . . .	143
<b>5</b>	<b>Huber Estimation</b>	<b>153</b>
5.1	Introduction . . . . .	153
5.2	Maximal Risk on $\mathbb{H}_d(\beta, L, M)$ . . . . .	158
5.3	Bandwidth Selector of Huber Estimator . . . . .	159
5.4	Proofs of Main Results: Upper Bounds . . . . .	161
5.4.1	Auxiliary Results: Large Deviations for M-estimators . . . . .	161
5.4.2	Proof of Proposition 7 . . . . .	164
5.4.3	Proof of Theorem 15 . . . . .	165
5.4.4	Proof of Theorem 16 . . . . .	166
5.5	Appendix . . . . .	169



# Chapitre 1

## Régression Non-Paramétrique

Nous présentons, dans ce chapitre introductif, les modèles de régression non-paramétriques que nous étudions. Nous introduirons l’approche localement paramétrique, le risque mini-max et la notion d’adaptation. Les estimateurs, utilisés dans cette thèse, sont présentés dans la section 1.3. Une partie importante de l’introduction est consacrée à une présentation détaillée de la méthode dite de *Lepski* (voir Section 1.5.2).

### 1.1 Objet de la Thèse

Dans cette thèse, nous étudions un domaine de la statistique mathématique : *l’estimation non-paramétrique*. Ceci consiste à estimer des fonctions (objets de dimensions infinies) à partir d’observations “bruitées”. Ce genre d’approche s’est considérablement développé ces dernières années dans le monde scientifique. L’imagerie (médicale ou astronomique), l’étude du génome (puces à ADN en grande dimension) ou encore les problèmes inverses (physique des matériaux, tomographie en imagerie médicale, etc.) ont recours à l’estimation non-paramétrique. Les statistiques ont un attrait particulier du fait qu’elles utilisent des théories mathématiques pour modéliser des problèmes réels. En particulier, l’introduction, de la notion d’*aléa* (en anglais : *random*) dans les observations, permet au statisticien de construire des méthodes d’estimation “fiables” en théorie.

La modélisation de problème se fait de la manière suivante. On dispose d’un nombre  $n$  de données, notées  $Y = (Y_1, Y_2, \dots, Y_n) \in \mathbb{R}^n$ , issues d’une expérience réelle. On modélise les observations en utilisant deux théories mathématiques, l’analyse fonctionnelle et les probabilités. On peut décomposer les observations  $Y$  de la façon suivante :

$$Y_i = f(X_i) + \text{“bruit”}, \quad i = 1, \dots, n, \quad X_i \in \mathbb{R}^d,$$

que l’on appelle *modèle additif*, où  $f$  est une fonction “régulière” à  $d$  variables de  $\mathbb{R}^d \rightarrow \mathbb{R}$ . Nous supposons pour toute cette thèse que la fonction  $f$  est dans un certain espace fonctionnel noté  $\mathcal{F}_\beta$  de dimension infinie où  $\beta$  est un paramètre de régularité.

On suppose également que nos observations sont dégradées par un aléa inconnu que nous modélisons par une variable aléatoire comme nous le voyons souvent en statistique (voir Section 1.2).

A partir de ce modèle mathématique, on veut résoudre le problème suivant : Estimer ou reconstruire la fonction  $f$  à partir des observations  $(Y_1, Y_2, \dots, Y_n)$ . Pour cela nous développons des outils mathématiques adéquats. Dans les modèles additifs, comme le modèle gaussien, des méthodes d'estimation ont été largement développées depuis un demi-siècle et sont devenus un outil courant de l'estimation non-paramétrique : estimateurs à noyau (voir Rosenblatt [1956], Parzen [1962], Nadaraya [1964], Watson [1964] et Borovkov [1987]) ou plus récemment la décomposition en bases d'ondelettes (Donoho, Johnstone, Kerkycharian, et Picard [1995] et Härdle, Kerkycharian, Picard, et Tsybakov [1998]).

Il est bien connu que les estimateurs linéaires (par rapport aux observations, voir Définition 4) ne sont pas *robustes* (non-sensibles aux valeurs extrêmes, en anglais *outliers*). Par exemple, si le bruit est une variable aléatoire de Cauchy (sans moment d'ordre 1), ces estimateurs sont inefficaces (Voir Exemple 1.1).

Si les observations admettent un moment d'ordre deux ( $\mathbb{E}Y^2 < \infty$ ), on sait, d'après le théorème central limite, que la moyenne empirique converge (en probabilité) vers son espérance à la vitesse  $1/\sqrt{n}$ . Ainsi, les estimateurs linéaires ne feront jamais mieux que cette vitesse. Ce qui implique que les estimateurs linéaires ne sont pas optimaux pour certains modèles (voir Exemple 1.1 et Section 1.2).

L'utilisation du phénomène de Stein [1981] repose sur le fait que le bruit est gaussien. En effet, Stein [1981] met au évidence la contraction de Stein dans le modèle de suites gaussiennes qui n'est valable seulement si le bruit est gaussien ou sous-gaussien. Bien que l'estimateur de James/Stein soit "meilleur" que la moyenne empirique, ces propriétés théoriques sont valables seulement dans le cas Gaussien.

L'adaptation, en estimation non-paramétrique, est aussi très "gourmande" en bruit gaussien, notamment pour l'obtention d'inégalités exponentielles de concentration. Par exemple pour l'estimateur à noyau, ceci est nécessaire (voir Section 1.5).

**Exemple 1.** Dans le cas paramétrique, on peut trouver des modèles dans lesquels les estimateurs linéaires ne sont pas optimaux. Par exemple, si les observations suivent une loi uniforme continue sur l'intervalle  $[0, \theta]$ ,  $Y_i \sim \mathcal{U}_{[0, \theta]}$ ,  $i = 1, \dots, n$ , alors la moyenne empirique notée  $\bar{Y}$  permet d'estimer  $\theta \geq 0$  avec la vitesse  $1/\sqrt{n}$ ,

$$2\bar{Y} \xrightarrow[n \rightarrow \infty]{n^{-1/2}} \theta, \quad \text{en probabilité.}$$

On peut construire un estimateur plus rapide dans ce modèle,

$$\max_i Y_i \xrightarrow[n \rightarrow \infty]{n^{-1}} \theta, \quad \text{en probabilité,}$$

un estimateur *non-linéaire* qui atteint la vitesse  $1/n$ . Nous étudions certains modèles avec cette particularité dans le cas non-paramétrique et nous développons un estimateur *bayésien*

non-linéaire. Celui-ci nous permet d'atteindre les vitesses de convergence optimales (Voir Définition de l'optimalité, Section 1.4).

**Exemple 2.** On regarde maintenant le modèle de Cauchy, pour lequel, beaucoup d'estimateurs ont échoué, notamment les estimateurs linéaires. On prend les variables  $Y_i = \theta + C_i$ ,  $i = 1, \dots, n$  avec  $C_i$  une variable de Cauchy (Définition 1.2.4). Il est facile de voir que la moyenne empirique  $\bar{Y}$  n'est pas convergente. Dans ce cas, on utilise la médiane empirique  $\text{Med}(Y)$  et sa normalité asymptotique (voir par exemple Brown, Cai, et Zhou [2008]). Ainsi, la médiane empirique converge vers  $\theta$  à la vitesse  $1/\sqrt{n}$ ,

$$\text{Med}(Y) \xrightarrow[n \rightarrow \infty]{n^{-1/2}} \theta, \quad \text{en probabilité.}$$

Pour le cas non-paramétrique, nous proposons un estimateur fondé sur cette idée que nous appelons *estimateur de Huber*.

Nous venons de présenter deux exemples où les estimateurs linéaires ne sont pas efficaces (voir inutilisables). Ces remarques motivent le travail de cette thèse avec l'introduction de nouveaux estimateurs non-linéaires.

Dans cette thèse, nous développons de nouveaux estimateurs localement paramétriques capables (sous certaines conditions) de s'adapter aux différents bruits (gaussien, Cauchy, etc.) et à la forme du modèle (additif ou multiplicatif). Ainsi, nous présentons deux types d'estimateurs : *estimateur bayésien* et *estimateur de Huber*. Nous montrons que pour différents modèles de régression, ces estimateurs sont optimaux au sens *minimax* (Définition 8) sur les *espaces de Hölder isotropes* (Définition 1).

La notion d'*adaptation*, introduite depuis une vingtaine d'années, est un point incontournable en estimation non-paramétrique (Stone [1982] et Efremovitch et Pinsker [1984]). Nous proposons plusieurs procédures adaptatives reposant sur la méthode dite de Lepski pour le choix de la fenêtre. Ces procédures permettent aux estimateurs considérés d'atteindre des vitesses de convergences adaptatives optimales en un certain sens (voir Section 4.3). L'utilisation de la méthode dite de Lepski est fréquente dans cette thèse. De ce fait, nous présentons dans ce chapitre, l'idée de la méthode ainsi que les conditions suffisantes à la mise en oeuvre de la procédure.

L'intérêt de l'estimation dépasse largement le cadre pratique et algorithmique. Il y a un vrai intérêt à développer des procédures qui sont optimales en théorie et de pouvoir les utiliser dans différents domaines d'applications. Ainsi, les performances théoriques de nos estimateurs adaptatifs reposent sur le contrôle des *grandes déviations*. Ce contrôle peut être obtenu par des inégalités de type *inégalités de concentration* pour les processus empiriques qui ont été très développées, entre autres par Talagrand [1995, 1996a, 1996b], Ledoux [1997], Birgé et Massart [1998], Massart [2000, 2007], Bousquet [2002], Boucheron, Bousquet, et Lugosi [2004], Golubev et Spokoiny [2009] et Goldenshluger et Lepski [2009b], du fait de leur nécessité.



## 1.2 Modèles de Régression et Espaces Fonctionnels

Dans cette section, nous présentons les espaces fonctionnels utilisés (espaces de Hölder), ainsi que les différents modèles de régression étudiés dans cette thèse. Dans un premier temps, nous parlerons du modèle de *régression générale*, et différents exemples de celui-ci, comme la régression additive avec densité inconnue, la régression gaussienne, la régression  $\alpha$  et la régression avec bruit multiplicatif uniforme. Un rapide survol est donné sur les modèles des statistiques non-paramétriques en fin de section.

### 1.2.1 Espaces de Hölder Isotropes

Tout au long de ce manuscrit, nous travaillerons exclusivement avec les espaces de Hölder isotropes, i.e. la fonction de régression inconnue est supposée à plusieurs variables avec la même régularité hölderienne dans chaque direction (par rapport à chaque variable).

Pour tout  $(p_1, \dots, p_d) \in \mathbb{N}^d$  nous notons  $\vec{p} = (p_1, \dots, p_d)$  et  $|\vec{p}| = p_1 + \dots + p_d$ .

**Définition 1.** Soient  $\beta > 0$ ,  $L > 0$ ,  $M > 0$ ,  $d \in \mathbb{N}^*$  et  $\lfloor \beta \rfloor$  le plus grand entier strictement inférieur à  $\beta$ . La classe de Hölder isotrope  $\mathbb{H}_d(\beta, L, M)$  est l'ensemble des fonctions  $f : [0, 1]^d \rightarrow \mathbb{R}$  ayant sur  $[0, 1]^d$  toutes ses dérivées d'ordre  $\lfloor \beta \rfloor$  et telles que  $\forall x, y \in [0, 1]^d$

$$\sum_{m=0}^{\lfloor \beta \rfloor} \sum_{|\vec{p}|=m} \sup_{x \in [0, 1]^d} \left| \frac{\partial^{|\vec{p}|} f(x)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} \right| \leq M,$$

$$\left| \frac{\partial^{|\vec{p}|} f(x)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} - \frac{\partial^{|\vec{p}|} f(y)}{\partial y_1^{p_1} \dots \partial y_d^{p_d}} \right| \leq L \|x - y\|_1^{\beta - \lfloor \beta \rfloor}, \quad \forall |\vec{p}| = \lfloor \beta \rfloor.$$

La théorie Minimax repose essentiellement sur le fait que la fonction à estimer est dans un espace fonctionnel, ici  $\mathbb{H}_d(\beta, L, M)$ . Les résultats minimax présentés dans la suite peuvent être étendus aux espaces anisotropes. Mais nous verrons dans la section 2.4 que les méthodes adaptatives ne sont pas conçues pour les espaces anisotropes dans le cas où l'estimateur est non-linéaire.

Pour étudier les modèles de régression inhomogène de Poisson et de régression multiplicative uniforme, une contrainte supplémentaire doit être rajoutée sur le support de  $f \in \mathbb{H}_d(\beta, L, M, A)$  où

$$\mathbb{H}_d(\beta, L, M, A) = \left\{ f \in \mathbb{H}_d(\beta, L, M) : \inf_{x \in [0, 1]^d} f(x) \geq A \right\}, \quad A > 0.$$

**Remarque 1.** Cette hypothèse est suffisante pour l'utilisation de l'estimateur bayésien pour les régressions inhomogène de Poisson et multiplicative uniforme, celle-ci est purement théorique (nécessaire) et n'est pas justifiée en pratique. Si on prend la régression multiplicative uniforme (définie dans la section 1.2.7), et si  $f$  admet des points de nullité, alors

$Y_i = 0$ . Avec cette observation, il est difficile de dire si cela est dû à la fonction ou au bruit (qui a une probabilité non-négligeable d'être proche de 0). Pour la régression inhomogène de Poisson, le paramètre de la loi de Poisson est toujours positif. Ainsi en pratique, on prendra  $A$  proche de 0 pour minimiser cette restriction. Sans cette restriction, le problème est ouvert.

La classe des fonctions de Hölder est incluse dans les espaces de Besov qui sont souvent utilisés dans l'approche maxiset ou avec les estimateurs par ondelettes (Autin [2004]). Nous donnons la définition de ces espaces ci-dessous.

**Les espaces de Besov Multidimensionnels.** Nous rappelons la définition des *espaces de Besov* pour les fonctions unidimensionnelles et quelques inclusions avec les espaces de Hölder et Sobolev. On définit les espaces de Besov, dans le cas de fonctions de  $[0, 1]^d$  dans  $\mathbb{R}$ . Il faut pour cela caractériser le *module de continuité* d'une fonction  $f$  de  $L^p([0, 1]^d)$ . Pour tout  $x$  dans  $[0, 1]^d$ , notons  $\Delta_h f(x) = f(x - h) - f(x)$ , et pour tout entier  $u$ , on note l'itérée  $\Delta_h^u f = \Delta_h \circ \dots \circ \Delta_h f$ . On définit alors le *module de continuité* pour la norme  $p$  (avec  $p \in [1, \infty]$ ) et pour tout  $t > 0$  de la manière suivante

$$\omega^p(f, t) = \sup_{\|h\|_2 \leq t} \left( \int_{J_{u,h}} |\Delta_h^u f(x)|^p dx \right)^{1/p},$$

où  $J_{u,h} = \{x \in [0, 1]^d, x + uh \in [0, 1]^d\}$  et  $\|\cdot\|_2$  est la norme  $\ell_2$  sur  $\mathbb{R}^d$ .

**Définition 2.** Soient  $p \in [1, \infty]$ ,  $q \in [1, \infty]$ ,  $s \in ]0, \infty[$  et  $u = \lceil s \rceil$  (où  $\lceil s \rceil$  est le plus petit entier strictement plus grand que  $s$ ). On dit qu'une fonction  $f$  appartenant à  $L^p([0, 1]^d)$  est dans l'espace de Besov  $\mathcal{B}_{p,q}^s([0, 1]^d)$ , quand  $\gamma_{spq}(f) < \infty$  où

$$\gamma_{spq}(f) = \begin{cases} \int_0^\infty (t^{-s} \omega^p(f, t))^q \frac{dt}{t}, & \text{si } 1 \leq q < \infty, \\ \sup_t |t^{-s} \omega^p(f, t)|, & \text{si } q = \infty. \end{cases}$$

On note  $\mathcal{B}_{p,q}^s(L)$  une boule de Besov de rayon  $L$ , munie de la norme :

$$\|f\|_{\mathcal{B}_{p,q}^s} = \|f\|_p + \gamma_{spq}(f).$$

Les espaces de Besov constituent une très grande famille de fonctions. En particulier, rappelons que l'espace de Sobolev  $S^s$  correspond précisément à l'espace  $\mathcal{B}_{2,2}^s$  et l'espace de Hölder  $\mathbb{H}_d(s, L)$  (avec  $0 < s \notin \mathbb{N}$ ) à l'espace  $\mathcal{B}_{\infty,\infty}^s(L)$  où  $L$  est le rayon de la boule de Besov. Ces espaces sont très utilisés dans l'approche par ondelettes. En effet, il est possible de définir les espaces de Besov à partir des coefficients dans la base d'ondelettes (voir Härdle, Kerkycharian, Picard, et Tsybakov [1998]). Pour plus de détails sur ces espaces, on se référera aux travaux de Bergh et Löfström [1976], Peetre [1976], Meyer [1992] ou DeVore et Lorentz [1993].

### 1.2.2 Régression Générale

Dans notre modèle, on observe les couples de variables aléatoires indépendantes  $(X_1, Y_1), \dots, (X_n, Y_n)$  notées

$$(1.2.1) \quad \mathcal{Z}_n = (X_i, Y_i)_{i=1, \dots, n},$$

où  $X_i$  est un vecteur dit de variables explicatives (appelé *design*) qui détermine la distribution de l'observation  $Y_i$ . Le vecteur  $X_i \in [0, 1]^d$  de dimension  $d$  peut être vu comme une variable temporelle ou spatiale et  $Y_i \in \mathbb{R}$  l'observation au point  $X_i$ . Notre modèle suppose que les valeurs  $X_i$  peuvent être aléatoires ou fixées et que la distribution de chaque  $Y_i$  est déterminée par un paramètre  $f_i$  qui peut dépendre de la position  $X_i$ ,  $f_i = f(X_i)$ . Dans beaucoup de cas, la paramétrisation naturelle est choisie de la façon suivante  $f_i = \mathbb{E}(Y_i|X_i)$  ( $\mathbb{E}$  est l'espérance mathématique). On note  $g(\cdot, f_i)$  la densité sur  $\mathbb{R}$  de l'observation  $Y_i$  par rapport à la mesure de Lebesgue. Pour la *régression additive*, nous considérerons le design aléatoire de loi uniforme sur  $[0, 1]^d$ . Le problème d'estimation est de reconstruire la fonction  $f$  en tout point  $y$ . Ce modèle que nous appelons *Régression générale*, sera traité dans le Chapitre 3.

Un estimateur dit *bayésien* est construit pour ce modèle très général. Nous verrons que sous certaines conditions (voir hypothèses 3), cet estimateur atteint les vitesses de convergence adaptatives optimales. Dans cette thèse, cette approche introduite par [Has'minskii et Ibragimov \[1981\]](#), pour l'estimation paramétrique, est généralisée à l'adaptation. Certains modèles, où la densité  $g$  des observations est discontinue, sont étudiés.

### 1.2.3 Régression Additive

Dans ce modèle, on observe les couples de variables aléatoires indépendantes  $(X_i, Y_i)_i$  qui vérifient l'équation :

$$(1.2.2) \quad Y_i = f(X_i) + \xi_i, \quad i = 1, \dots, n,$$

où le bruit est une variable aléatoire de densité  $g_\xi(\cdot)$ .  $X_i$  est le design aléatoire de loi uniforme sur  $[0, 1]^d$  et  $X_i$  est indépendant de  $\xi_i$ . La fonction de régression  $f$  est supposée appartenir à une boule de Hölder  $\mathbb{H}_d(\beta, L, M)$ . Nous supposerons vérifiées les hypothèses suivantes :

#### Hypothèses 1.

1.  $g_\xi$  est paire,
2.  $g_\xi(0) \geq A > 0$ ,
3.  $g_\xi$  est continue en 0.

**Remarque 2.** La symétrie du bruit (hypothèse 1.1) est une supposition assez faible. En général, on suppose que les effets du bruit sont les mêmes à gauche et à droite de la médiane qui vaut 0. Notons qu'un grand nombre de densités, utilisées dans cette thèse, vérifient cette

*hypothèse. Les deux autres hypothèses sont nécessaires pour contrôler les grandes déviations de l'estimateur de Huber (voir Chapitre 5). Elles permettent aussi de vérifier l'assertion suivante. Avec ces hypothèses, on peut voir que  $\xi$  admet une médiane théorique unique égale à 0, i.e. 0 est unique solution de l'équation  $\mathbb{P}(\xi < x) = 1/2$ . Noter qu'il existe de nombreux modèles classiques vérifiant ces hypothèses (Régressions gaussienne, de Cauchy et  $\alpha$ ). On peut remarquer que la régression additive est un cas particulier de la régression générale. En effet, il suffit de prendre  $g(\cdot, f(X_i)) = g_\xi(\cdot - f(X_i))$ .*

L'étude de ce modèle nécessite la construction d'un estimateur fondé sur l'idée de la médiane (estimateur de Huber, voir Chapitre 5). En particulier, nous développons une procédure adaptative dans le cas où la régularité de la fonction cible est inconnue.

Hall et Jones [1990] ont utilisé une méthode de *Validation croisée* (voir Tsybakov [2008]) sur une famille d'estimateurs robustes pour obtenir des résultats adaptatifs avec le risque  $L_2$ . Peu après, Härdle et Tsybakov [1992] ont étendu ce résultat, avec une méthode de *Plug-in* (introduit par Woodroffe [1970]), pour des *fonctions de contrast* plus générales, avec un choix aléatoire local de la fenêtre, mais seulement des résultats de normalité asymptotique sont donnés.

L'adaptation dans la régression additive a fait l'objet d'un autre travail. En effet, Brown, Cai, et Zhou [2008] utilisent la normalité asymptotique de la médiane pour approximer ce modèle par le modèle gaussien avec une méthode de médiane par blocs. Une étape intermédiaire est de projeter les nouvelles observations dans une base d'ondelettes. Ensuite, la méthode de Stein par blocs (voir Cai [1999]) est utilisée pour l'adaptation en estimation globale. En revanche, cette approche nécessite des hypothèses plus fortes que les hypothèses 1 qui sont suffisantes pour l'utilisation de l'estimateur de Huber.

Plus récemment, Reiss, Rozenholc, et Cuenod [2009] utilise la méthode de Lepski pour développer un estimateur de Huber adaptatif dans le modèle additif. Les résultats obtenus sont pour l'estimation ponctuelle mais pour les fonctions localement constante (i.e.  $\beta \leq 1$ ).

Dans la suite, nous présentons plusieurs cas particuliers où les estimateurs bayésien et de Huber peuvent être utilisés. En effet nous vérifions dans le cas de la régression gaussienne, si nos estimateurs sont applicables. Pour le bruit de Cauchy, l'estimateur de Huber est tout désigné pour estimer la fonction de régression (voir Chapitre 5). Nous présentons des modèles peu connus dans la littérature (voir Régressions  $\alpha$  et multiplicative uniforme). Nous utilisons l'estimateur *bayésien* dans le cadre du modèle inhomogène de Poisson (très utilisé en imagerie, voir notamment Polzehl et Spokoiny [2006] et Katkovnik et Spokoiny [2008]). Ici la vitesse de convergence est la même que pour le modèle gaussien (qui peut être obtenu avec un estimateur linéaire). Les régressions, dites  $\alpha$  et *multiplicative uniforme*, sont de très bons exemples dans lesquels la vitesse de convergence devient meilleure que la vitesse des estimateurs linéaires. Pour plus de détails, voir Section 2.1.4.

### 1.2.4 Régression Gaussienne et de Cauchy

Les régressions *gaussienne* ou de *Cauchy* sont des cas particuliers de la *régression additive* où la densité  $g_\xi(\cdot)$  prend respectivement les formes suivantes

$$\text{Densité gaussienne : } g_\xi(x) = (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{x^2}{2\sigma^2}\right\}, \quad \sigma > 0,$$

$$\text{Densité de Cauchy : } g_\xi(x) = \frac{1}{\pi} \frac{a}{a^2 + x^2}, \quad a > 0.$$

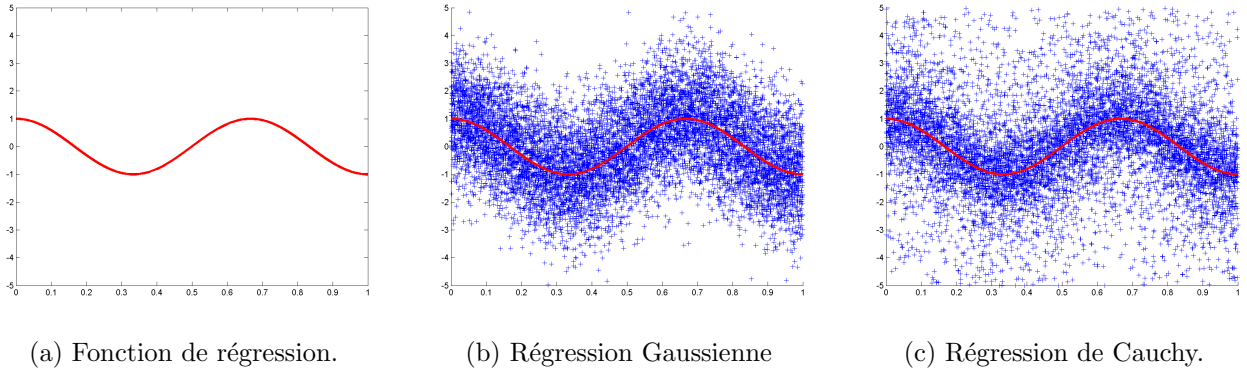


FIGURE 1.1 – Observations dans les régressions gaussiennes et de Cauchy.

Néanmoins, elles suscitent un intérêt particulier. En effet le bruit gaussien est utilisé de façon systématique dans le domaine des statistiques et le bruit de Cauchy est une variable aléatoire sans moment d'ordre 1. Pour le bruit gaussien, on peut trouver une multitude d'articles qui l'étudient. Citons deux livres de [Nemirovski \[2000\]](#) et [Tsybakov \[2008\]](#) pour une introduction à ce modèle et aux méthodes classiques d'estimation non-paramétrique. Nous introduisons ce modèle standard pour vérifier que les estimateurs développés dans cette thèse fonctionnent correctement dans ce modèle classique (voir Chapitres 3 et 5).

L'estimateur de *Huber* est particulièrement bien adapté aux bruits dont les densités sont à queues lourdes (par exemple bruit de Cauchy, voir Chapitre 5). On peut constater une différence notable, entre les bruits gaussien et de Cauchy, dans la figure 1.1. En effet pour le bruit de Cauchy, on constate un nombre plus important de *valeurs extrêmes*. Ceci explique l'inefficacité des estimateurs linéaires pour ce bruit.

### 1.2.5 Régression Inhomogène de Poisson

Considérons la régression générale, cette fois avec la particularité que les observations sont discrètes  $Y_i \in \mathbb{N}$ . Nous supposons que  $Y_i$  suit une loi de Poisson de paramètre  $f(X_i)$

$(Y_i \sim \mathcal{P}(f(X_i)))$ . On écrit la densité

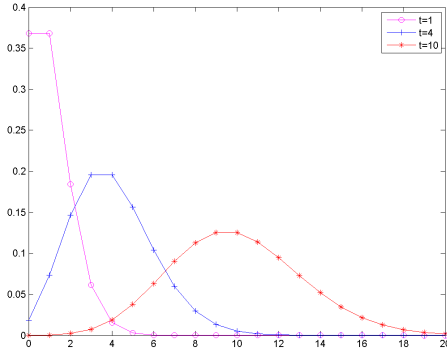
$$g(k, f(X_i)) = \mathbb{P}_f(Y_i = k) = \frac{[f(X_i)]^k}{k!} \exp\{-f(X_i)\}, \quad k \in \mathbb{N}, \quad f \in \mathbb{H}_d(\beta, L, M, A).$$

Les points du design  $(X_i)_{i \in 1, \dots, n}$  sont déterministes et sans perte de généralité nous supposons qu'ils sont répartis de manière uniforme sur la grille suivante :

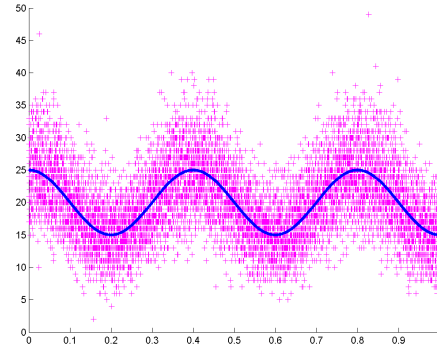
$$(1.2.3) \quad X_i \in \{1/n^{1/d}, 2/n^{1/d}, \dots, 1\}^d, \quad i = 1, \dots, n.$$

Remarquons la restriction de  $f$  à l'espace  $\mathbb{H}_d(\beta, L, M, A)$ , cette restriction ( $f \geq A$ ) est retrouvée dans la régression multiplicative uniforme. Ce modèle est très utilisé en imagerie, en particulier pour modéliser la photométrie des appareils photo numérique. On le trouve aussi en tomographie, il modélise la résonance magnétique des positons utilisée pour obtenir des IRM d'une ou plusieurs parties du corps humain.

Les travaux de [Anscombe \[1948\]](#) permettent de passer approximativement de la régression de Poisson à la gaussienne par une transformation des observations par la fonction  $x \rightarrow 2\sqrt{x + 3/8}$ . Dans la figure 1.2, nous donnons des exemples de densités de Poisson et un échantillon d'observations simulé à partir d'une fonction de régression.



(a) Densités d'une loi de Poisson de paramètre  $t$ .



(b) Observations de Poisson.

FIGURE 1.2 – Densité de la loi de Poisson de paramètre  $t > 0$  et observations.

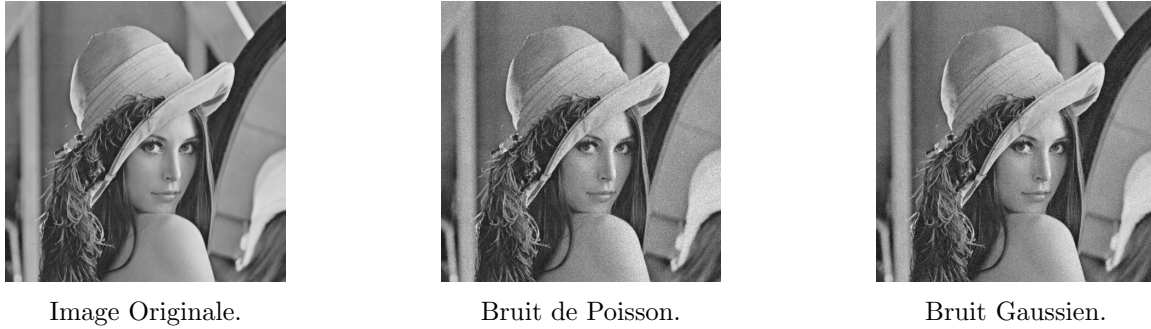


FIGURE 1.3 – Illustration d’images issues d’observations de Poisson (zoomer sur la version numérique).

### 1.2.6 Régression $\alpha$

Le modèle de régression  $\alpha$  est le suivant

$$(1.2.4) \quad Y_i = f(X_i) + \epsilon_i, \quad i = 1, \dots, n,$$

où  $\epsilon_i$  a pour densité  $g_\alpha(x) = C(\alpha) \exp \{-|x|^\alpha\}$  avec  $0 < \alpha < 1/2$ .  $C(\alpha)$  est une constante choisie pour que  $g_\alpha$  soit bien une densité. Ici le design  $X_i$  est choisi comme dans (1.2.3) et  $f \in \mathbb{H}_d(\beta, L, M)$ . La particularité de ce modèle est que pour  $\alpha < 1/2$  la vitesse de convergence change en fonction de  $\alpha$  (voir Chapitre 3).

La figure 1.4 présente les observations de ce modèle pour  $\alpha = 1/4$ . On constate qu’il y a un très grand nombre d’observations éloignées de la fonction de régression  $f$ . Mais un petit nombre est très proche de  $f$  (à comparer avec la figure 1.1), cela nous donne une information très précieuse qui peut se vérifier en théorie à travers de meilleures vitesses de convergence. En effet, nous verrons que l’estimateur bayésien est plus rapide (au sens des vitesses de convergence) que les estimateurs linéaires (voir Sections 2.1 et 3.4.3).

Nous donnons la figure 1.5 pour illustrer ce comportement. En effet, la queue de distribution ne joue pas un rôle majeur entre les densités. Le comportement de la densité, à mettre en avant, est la concentration du bruit au voisinage de 0. Le voisinage de la gaussienne (on parle de l’intervalle  $[-3, 3]$ ) est assez large, tandis que le voisinage de  $g_\alpha$  est plus concentré autour de 0 (on constate un pic en 0 de la fonction  $g_\alpha$ ). Donc bien que la queue de distribution gaussienne soit la plus faible, le comportement autour de 0, handicape cette loi. Nous avons aussi fait apparaître la loi de Cauchy qui a le même comportement que la gaussienne au voisinage de 0, mais une queue de distribution lourde.

Dans la figure 1.5, nous avons choisi les paramètres de chaque densité, pour que toutes soient égales à 1 au point 0. Ceci est utilisé comme critère de comparaison de la dispersion (ou la variance, mais n’existe pas pour la loi de Cauchy) de chaque densité.



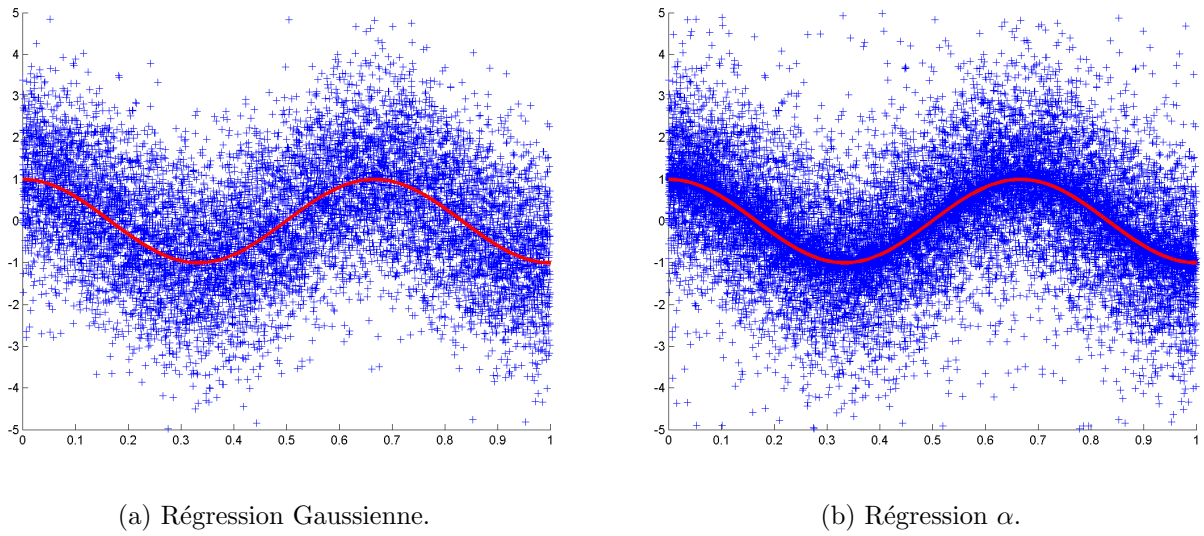


FIGURE 1.4 – Comparaison entre des observations gaussiennes et des observations issues de la régression  $\alpha$  pour  $\alpha = 1/4$ .

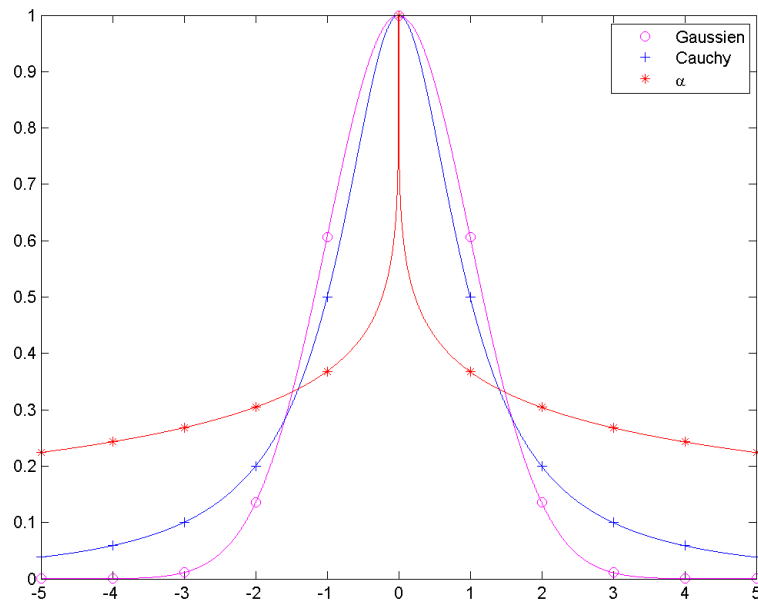


FIGURE 1.5 – Densités gaussienne, de Cauchy et  $\alpha$ ,  $\alpha = 1/4$



### 1.2.7 Régression Multiplicative Uniforme

Soit l'expérience statistique générée par les observations  $\mathcal{Z}_n = (X_i, Y_i)_{i=1, \dots, n}$ ,  $n \in \mathbb{N}^*$ , où  $(X_i, Y_i)$  satisfont à l'équation

$$(1.2.5) \quad Y_i = f(X_i) \times U_i, \quad i = 1, \dots, n,$$

où  $f \in \mathbb{H}_d(\beta, L, M, A)$ . Les variables aléatoires  $(U_i)_{i \in 1, \dots, n}$  sont supposées indépendantes et uniformément distribuées sur  $[0, 1]$ . Le design  $X_i$  est choisi comme dans (1.2.3). Ce modèle est étudié dans le chapitre 4 pour lequel on utilise l'estimateur *bayésien* pour estimer la fonction de régression. Ce type de modèle peut être utilisé dans le cadre des modèles de frontière développés par Simar et Wilson [2000]. Outre le cadre pratique, il est intéressant de voir que les estimateurs linéaires sont inefficaces dans ce modèle (voir Section 2).

Nous verrons que pour ce modèle, la vitesse de convergence minimax est très rapide. Elle est meilleure que toutes celles des modèles précédents. Visuellement, on peut le constater à l'aide de la figure 1.6. En effet, toutes les observations sont en dessous de la fonction de régression  $f$ . Et on voit distinctement apparaître la fonction cible, en regardant le maximum par morceaux des observations. L'estimateur bayésien sera, dans ce modèle, beaucoup plus rapide que les estimateurs linéaires (Voir Sections 2.1 et 3.4.4).

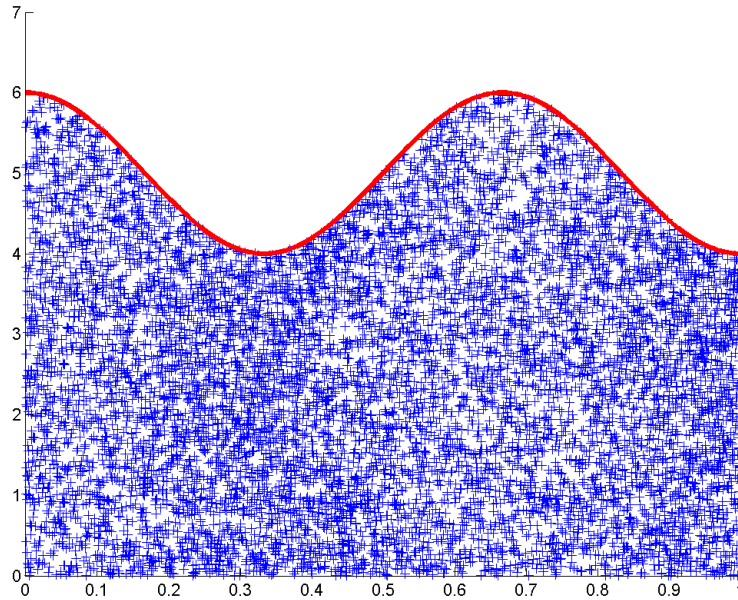


FIGURE 1.6 – Observations d'une fonction bruitée avec un bruit multiplicatif uniforme.

**Autres Modèles.** On peut citer d'autres modèles des statistiques non-paramétriques : le modèle du bruit blanc et les problèmes inverses. Le modèle de bruit blanc gaussien est défini par l'équation différentielle stochastique :

$$dY_t = f(t)dt + \frac{\sigma}{\sqrt{n}}dW_t, \quad t \in \mathbb{R}^d,$$

où  $f$  est la fonction à estimer à partir des observations  $(Y_t)_{t \in \mathbb{R}^d}$ ,  $W_t$  est un champ Brownien sur  $\mathbb{R}^d$  (Processus de Wiener),  $\sigma > 0$  et  $n \in \mathbb{N}^*$ . Le modèle de bruit blanc gaussien joue un rôle important en Statistiques (voir [Has'minskii et Ibragimov \[1981\]](#)). L'avantage de ce modèle est qu'il est simple à utiliser et qu'il approche d'autres modèles statistiques, en particulier la régression *gaussienne* ([Brown et Low \[1996\]](#) et [Nussbaum \[1996\]](#)).

Pour ce modèle, les estimateurs à noyau sont utilisés pour l'adaptation (voir les travaux les plus récents de [Kerkycharian, Lepski, et Picard \[2001\]](#), [Klutchnikoff \[2005\]](#), [Goldenshluger et Lepski \[2008\]](#), [Goldenshluger et Lepski \[2009a\]](#)).

Les problèmes inverses sont le résultat de l'équation suivante :

$$dY_t = Af(t)dt + \frac{\sigma}{\sqrt{n}}dW_t, \quad t \in \mathbb{R}^d,$$

où  $A$  est un opérateur linéaire défini sur un espace fonctionnel. L'heuristique de ce modèle est que l'on observe indirectement la fonction cible, à travers une transformation de celle-ci, par l'opérateur  $A$  connu mais pas forcément inversible. Les problèmes inverses sont très connus pour leurs applications en pratique (équation de la chaleur, tomographie aux rayons X, etc.). Depuis ces dix dernières années, de nombreuses méthodes ont été développées pour l'estimation adaptative de la fonction  $f$  (voir par exemple [Cavalier et Tsybakov \[2002\]](#) et [Cavalier, Golubev, Picard, et Tsybakov \[2002\]](#)). Par exemple, la méthode de *Stein par blocs* est fréquemment utilisée. Le modèle est projeté dans une base et on obtient le modèle de suite gaussienne. On observe alors les coefficients de la fonction  $f$ , multipliés par les valeurs singulières de l'opérateur  $A$ , auxquels on ajoute un bruit gaussien. Une autre méthode, développée pour les problèmes inverses très mal posés (l'opérateur  $A$  n'est pas inversible), est appelée *enveloppe du risque* et introduite par [Cavalier et Golubev \[2006\]](#), [Cavalier \[2008\]](#) et [Marteau \[2009\]](#).

### 1.3 Approche Localement Paramétrique

Dans cette section, nous présentons l'approche *localement paramétrique* (en anglais : *local parametric fitting*). Plus précisément, nous traitons l'approche *polynomiale locale* qui nous permet d'approximer localement la fonction  $f$  par un polynôme noté  $f_\theta$  lui-même inconnu. Ensuite, nous proposons deux estimateurs permettant d'estimer les coefficients du polynôme d'approximation  $f_\theta$ . Cette technique est appelée *Méthode des polynômes locaux*, elle a été utilisée pour la première fois par [Katkovnik \[1985\]](#) (pour plus de détails voir [Tsybakov \[2008\]](#)). Notons que, en pratique, le degré du polynôme d'approximation est choisi assez petit (degré = 0, 1 ou encore 2), ceci est dû au temps de calculs qui est exponentiellement grand par rapport au nombre de coefficients à estimer.

L'approche localement paramétrique ne s'arrête pas qu'aux polynômes. On peut imaginer approximer la fonction  $f$  par tout autre objet paramétrique. Par exemple, on peut décomposer  $f$  dans une base orthogonale où l'on ne gardera qu'un nombre fini de coefficients.

**Méthode des Polynômes Locaux.** On se place dans le modèle de régression générale défini dans (1.2.1). Soit  $f \in \mathbb{H}_d(\beta, L, M)$  avec  $L, M > 0$  et  $\beta \in [0, b[$  où  $b \in \mathbb{N}^*$  peut être choisi arbitrairement grand ( $b$  est la régularité maximale des fonctions que l'on étudie). Soit

$$(1.3.1) \quad V_h(y) = \left\{ X_i \in \bigotimes_{j=1}^d [y_j - h/2, y_j + h/2] \cap [0, 1]^d \right\},$$

un voisinage de  $y$  tel que  $V_h(y) \subseteq [0, 1]^d$ , où est  $h \in (0, 1)$  un scalaire donné. Soit

$$(1.3.2) \quad D_b = \sum_{m=0}^b \binom{m+d-1}{d-1}.$$

Soit  $K(z)$ ,  $\forall z \in \mathbb{R}^d$  le vecteur de dimension  $D_b$  des polynômes du type suivant :

$$K^\top(z) = \left( \prod_{j=1}^d z_j^{p_j}, (p_1, \dots, p_d) \in \mathbb{N}^d : 0 \leq p_1 + \dots + p_d \leq b \right),$$

où le signe  $\top$  représente la fonction transposition. Ensuite, pour tout  $t \in \mathbb{R}^{D_b}$ , où  $t^\top = (t_{p_1, \dots, p_d}, (p_1, \dots, p_d) \in \mathbb{N}^d : 0 \leq p_1 + \dots + p_d \leq b)$ , on définit le polynôme local

$$(1.3.3) \quad f_t(x) = t^\top K \left( \frac{x-y}{h} \right) \mathbb{I}_{V_h(y)}(x), \quad x \in [0, 1]^d,$$

où  $\mathbb{I}$  est la fonction indicatrice. Notons que  $f_t(y) = t_{0, \dots, 0}$ . Introduisons l'ensemble des coefficients

$$(1.3.4) \quad \Theta(M) = \{t \in \mathbb{R}^{D_b} : \|t\|_1 \leq M\} \subset \mathbb{R}^{D_b},$$

où  $\|\cdot\|_1$  est la norme  $\ell_1$  sur  $\mathbb{R}^{D_b}$ . Remarquons que pour tout  $f \in \mathbb{H}_d(\beta, L, M)$

$$\exists \theta = \theta(f, y, h) \in \Theta(M) : \sup_{x \in V_h(y)} |f(x) - f_\theta(x)| \leq Ldh^\beta.$$

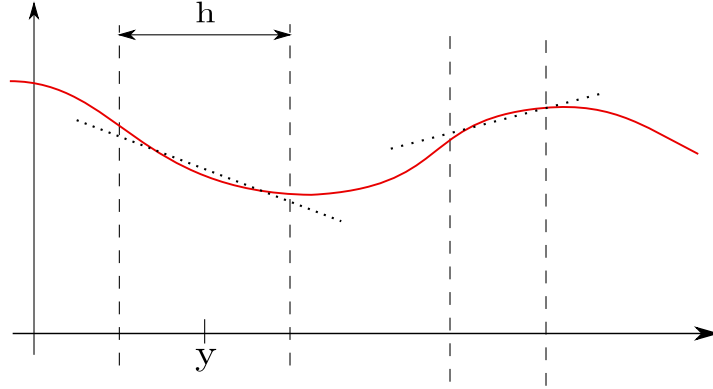


FIGURE 1.7 – Illustration d'une approximation locale d'une fonction (ligne) par une droite ou un polynôme d'ordre 1 (pointillé).

Le polynôme local  $f_\theta$  peut être vu comme une approximation localement paramétrique de la fonction de régression  $f$  sur le voisinage  $V_y(h)$ . Par exemple, on peut choisir pour  $f_\theta$  le polynôme de Taylor de la fonction  $f$  au point  $y$  (voir Figure 1.7).

Comme nous l'avons précisé en début de chapitre, la restriction aux polynômes locaux n'est pas nécessaire, si l'on trouve un autre objet paramétrique capable d'approcher la fonction cible (par exemple, une décomposition en base d'ondelettes à coefficients finis).

L'idée principale est que si  $h$  est choisi suffisamment petit, les observations originales (1.2.1) sont bien approximées par le modèle paramétrique  $\mathcal{Y}_i$  de densité  $g(\cdot, f_\theta(X_i))$  dans lequel, sous certaines conditions sur  $g(\cdot)$ , les estimateurs de Huber et bayésien sont optimaux au sens des vitesses de convergences. Dans la suite nous présentons les estimateurs qui nous permettent d'estimer  $\theta$ , en rappelant pour commencer, la définition d'un estimateur.

**Définition 3.** On dit que  $\tilde{f}$  est un estimateur si  $\tilde{f}(\cdot) = \tilde{f}(\cdot, \mathcal{Z}_n)$  est une fonction mesurable des observations.

**Définition 4.** On dit que  $\tilde{f}$  est un estimateur linéaire si il existe  $K : [0, 1]^d \times [0, 1]^d \rightarrow \mathbb{R}$  telle que :

$$\tilde{f}(y) = \sum_{i=1}^d K(X_i, y) Y_i, \quad \forall f \in \mathbb{H}_d(\beta, L), \quad \forall y \in [0, 1]^d.$$

On peut trouver cette définition dans les notes de [Nemirovski \[2000\]](#).

**Définition 5.** On appelle  $f_{\hat{\theta}}$  estimateur par polynômes locaux si  $\hat{\theta}$  est solution du problème de minimisation suivant :

$$\min_{t \in \Theta(M)} \ell(\mathcal{Z}_n, f_t),$$

où  $\ell(.,.)$  est un critère à choisir.

Dans cette thèse, on rencontre les critères bayésien, des moindres carrés, des valeurs absolues et de Huber (voir ci-dessous).

### 1.3.1 Estimateur Bayésien

Considérons, dans cette section, le modèle de régression *générale* définie dans la section 1.2.2 et soit  $\mathbb{E}_f = \mathbb{E}_f^n$  l'espérance mathématique par rapport à la loi de probabilité  $\mathbb{P}_f = \mathbb{P}_f^n$  des observations  $\mathcal{Z}_n$ . Définissons la *pseudo vraisemblance*

$$(1.3.5) \quad L_h(t, \mathcal{Z}_n) = \prod_{X_i \in V_h(y)} g(Y_i, f_t(X_i)), \quad t \in \Theta(M),$$

et le *critère bayésien*

$$(1.3.6) \quad \pi_h(t) = \int_{\Theta(M)} \|t - u\|_1 [L_h(u, \mathcal{Z}_n)]^{1/m} du, \quad t \in \Theta(M).$$

où  $m$  est une constante positive choisie dans la suite (voir Hypothèses 1, Chapitre 3).

Soit  $\hat{\theta}(h)$  la solution du problème de minimisation suivant :

$$(1.3.7) \quad \hat{\theta}(h) = \arg \min_{t \in \Theta(M)} \pi_h(t),$$

où  $\|\cdot\|_1$  est la norme  $\ell_1$  sur  $\mathbb{R}^{D_b}$ . L'estimateur localement bayésien  $\hat{f}^h(y)$  de  $f(y)$  est défini maintenant comme  $\hat{f}^h(y) = \hat{\theta}_{0,\dots,0}(h)$ .

Nous constatons que la vraisemblance joue le rôle de la loi *a priori* sur le paramètre à estimer. Il est possible de rajouter un autre *a priori* sur le paramètre. En effet, on a l'habitude de voir (Has'minskii et Ibragimov [1981]) un critère de la forme  $\int \|t - u\|_1 L_h(u) q(u) du$  où  $q(\cdot)$  est une densité sur  $\Theta(M)$  (dans notre cas  $q(\cdot) \equiv 1$ ). Dans l'approche bayésienne, le but est de choisir de façon optimale cette densité. Ici, nous nous intéressons seulement au problème d'optimalité des vitesses de convergence. On remarque aussi que la vraisemblance est élevée à la puissance  $1/m$  ( $m$  est à choisir). C'est une nouvelle façon de considérer l'estimateur de type *bayésien*, traditionnellement  $m = 1$ . Le choix de  $m$  permet d'affaiblir les hypothèses sur la densité  $g(\cdot)$  et d'avoir une meilleure compréhension de la démonstration. Dans le chapitre 4 nous prenons  $m = 1$ , en revanche dans le chapitre 3 nous considérons  $m$  à choisir (par exemple  $m = 2$  pour la régression gaussienne). Le chapitre 4 est un cas particulier du chapitre 3 du point de vue des modèles, mais les démonstrations sont plus complexes dans le chapitre 4 car  $m = 1$  est mal ajusté.

Le choix de la fonction de perte, ici la norme  $\ell_1$ , n'est pas restrictif. Ceci nous permet d'améliorer les constantes de majoration du risque. On peut prendre notamment la norme  $\ell_2$  qui permet de donner une forme explicite de  $\hat{\theta}$  pour l'estimation paramétrique unidimensionnelle ([Has'minskii et Ibragimov \[1981\]](#)) :

$$\hat{\theta}(h) = \frac{\int_{\Theta(M)} u [L_h(u, \mathcal{Z}_n)]^{1/m} du}{\int_{\Theta(M)} [L_h(v, \mathcal{Z}_n)]^{1/m} dv}.$$

Cet estimateur est aussi connu sous le nom d'*estimateur de Pitman* (voir [Has'minskii et Ibragimov \[1981\]](#)). Cette version de l'estimateur bayésien donne un estimateur linéaire pour la régression gaussienne. Mais cela ne dépend pas de la norme  $\ell_2$  utilisée ici. En effet la forme de l'estimateur bayésien dépend directement de la densité  $g(.,.)$  des observations et non de la fonction de perte utilisée pour sa construction. Notons que cette méthode est similaire à une approche localement paramétrique établie à partir d'estimateurs du maximum de vraisemblance récemment développés par [Polzehl et Spokoiny \[2006\]](#) et [Katkovnik et Spokoiny \[2008\]](#) pour les *modèles statistiques réguliers*.

Le principal avantage de ce nouvel estimateur *bayésien*, est qu'il atteint la vitesse de convergence optimale du modèle. Par exemple, pour les régressions *gaussienne* et *multiplicative uniforme*, notre estimateur atteint respectivement les vitesses minimax (voir Définition 8)  $n^{-\frac{\beta}{2\beta+d}}$  et  $n^{-\frac{\beta}{\beta+d}}$  sous les hypothèses 3. Nous pouvons dire que l'estimateur *bayésien* concurrence (fait mieux ou aussi bien que) les estimateurs linéaires dans certains modèles (voir Section 2.1).

### Estimateur du Maximum de Vraisemblance

Il est très important de se poser la question suivante : existe-t-il d'autres estimateurs optimaux comme l'estimateur bayésien ?

La réponse à cette question est positive. En effet, l'exemple le plus connu est l'estimateur du maximum de vraisemblance. Cet estimateur concurrence les estimateurs linéaires au même titre que l'estimateur bayésien. On peut expliquer ce phénomène par le fait que ces deux estimateurs utilisent l'information contenue dans la vraisemblance du modèle et donc sont capables d'améliorer la vitesse. Définissons cet estimateur

$$\hat{f}^h(y) = \arg \max_{t \in \Theta(M)} L_h(t, \mathcal{Z}_n)$$

Dans le cas paramétrique, si l'on considère le modèle  $Y_i \sim \mathcal{U}_{[0,\theta]}$ ,  $i = 1, \dots, n$ , l'estimateur du maximum de vraisemblance correspondant est  $\max_{i=1, \dots, n} Y_i$ . Comme cet estimateur est basé sur le maximum des observations, donc non-linéaire, il est possible d'améliorer la vitesse de convergence. Ainsi, l'estimateur  $\max_{i=1, \dots, n} Y_i$  converge vers  $\theta$  à la vitesse  $1/n$  contre  $1/\sqrt{n}$  pour les estimateurs linéaires.

On peut généraliser cette approche au cas non-paramétrique en utilisant comme estimateur  $\max_{X_i \in V_h(y)} Y_i$  (estimateur localement constant). Il semble très difficile de pouvoir généraliser cette approche pour un estimateur localement polynômial d'ordre  $D_b$  quelconque.

Un autre problème est qu'il faut supposer une hypothèse supplémentaire (par rapport à l'estimateur bayésien) de continuité sur la vraisemblance. [Has'minskii et Ibragimov \[1981\]](#) (voir Chapitre 1, Section 5, Théorème 5.1) supposent cette hypothèse pour l'estimation paramétrique. Mais jusqu'à présent, on ne trouve pas dans la littérature une extension de l'estimateur du maximum de vraisemblance local comme nous l'avons fait pour l'estimateur bayésien dans le problème non-paramétrique.

[Polzehl et Spokoiny \[2006\]](#) utilisent l'estimateur du maximum de vraisemblance dans les *modèles de famille exponentielles* dans le cadre localement paramétrique (fonctions localement constantes). Les auteurs donnent des résultats très généraux sur ces modèles en supposant l'existence de *l'information de Fisher (modèle régulier)*. Mais les vitesses de convergences ne sont pas améliorées et restent les mêmes que pour les estimateurs linéaires.

Ainsi, pour la régression multiplicative uniforme non-paramétrique, il n'est pas prouvé que l'on peut utiliser l'estimateur du maximum de vraisemblance et cela semble un objectif difficile.

Comme on vient de le voir, les raisons de l'utilisation de l'estimateur bayésien (préféré à l'estimateur du maximum de vraisemblance) sont essentiellement techniques.

### 1.3.2 Estimateur de Huber

On se place, dans cette section, dans le cadre de la régression additive avec la densité  $g_\xi$  inconnue définie dans (1.2.2). Nous supposons vraies les hypothèses 1 sur la densité  $g_\xi$ . Nous introduisons une variante de l'estimateur de la médiane, plus communément appelé *estimateur de Huber*, développé pour l'estimation non-paramétrique par [Tsybakov \[1982a, 1982b, 1983, 1986\]](#) et [Härdle et Tsybakov \[1988, 1992\]](#). Ce dernier donne une règle de sélection locale pour la fenêtre basée sur la méthode *Plug-in* avec des résultats en normalité asymptotique. Un peu plus tôt [Hall et Jones \[1990\]](#) propose d'utiliser la *validation croisée* pour le risque  $L_2$ . Plus récemment, [Reiss, Rozenholc, et Cuenod \[2009\]](#) ont développé l'adaptation d'estimateurs robustes, en particulier le critère de Huber mais seulement pour les estimateurs localement constants. Considérons la *fonction de Huber*,

$$(1.3.8) \quad Q(z) = \frac{z^2}{2} \mathbb{I}_{|z| \leq 1} + \left( |z| - \frac{1}{2} \right) \mathbb{I}_{|z| > 1},$$

où  $\mathbb{I}$  est la fonction indicatrice.

Définissons le *critère de Huber*

$$(1.3.9) \quad \tilde{m}_h(t) = \tilde{m}_h(t, \mathcal{Z}_n) := \frac{1}{nh^d} \sum_{i=1}^n Q(Y_i - f_t(X_i)) \mathbb{I}_{\{X_i \in V_h(y)\}}$$



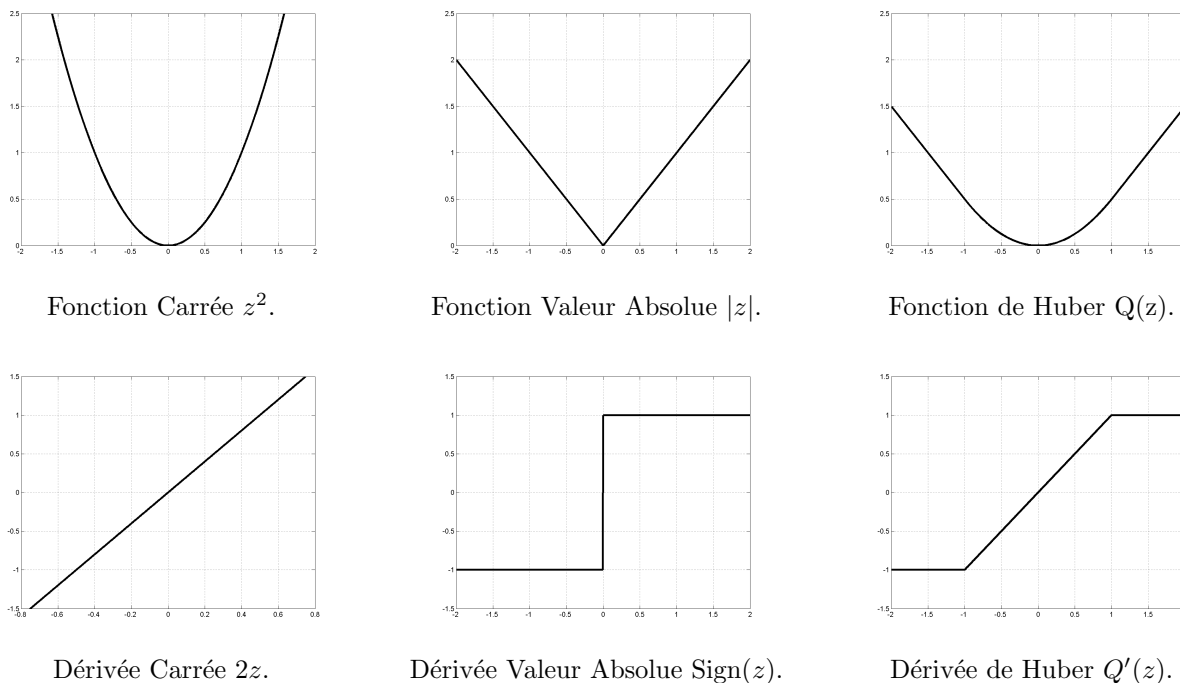


FIGURE 1.8 – Courbes des fonctions carré, valeur absolue, de Huber et leurs dérivées.

On introduit le critère de Huber qui allie les propriétés pratiques (robustesse) du critère  $\ell_1$  et les propriétés théoriques du critère  $\ell_2$  (dériverabilité). On parle d'estimateur *robuste* si l'estimateur en question n'est pas sensible aux valeurs extrêmes, (voir [Rousseeuw et Leroy \[1987\]](#) et [Huber et Ronchetti \[2009\]](#)).

Ce critère fait partie de la famille des critères par *polynômes locaux*. Par exemple, on peut utiliser le critère  $\ell_2$  appelée *moindres carrés* ( $Q(z) = z^2$ , voir [Tsybakov \[2008\]](#), Chapitre 1), ou bien le critère  $\ell_1$  appelé *valeurs absolues* ( $Q(z) = |z|$ ).

Le critère  $\ell_2$ , qui donne un estimateur fondé sur la moyenne, est beaucoup plus sensible aux valeurs extrêmes (voir Figure 1.8, la dérivée de la fonction carrée). En pratique, quand on cherche un estimateur robuste, on utilise le critère  $\ell_1$  qui conduit à un estimateur fondé sur la médiane. Mais en estimation non-paramétrique, cela pose un problème car la fonction valeur absolue n'est pas dérivable en 0 (voir Figure 1.8). L'avantage de la fonction de Huber est que sa dérivée est continue sur  $\mathbb{R}$ , tandis que la dérivée de la fonction valeur absolue n'est pas continue en 0 (voir Figure 1.8). On peut voir dans le chapitre 5, que le contrôle des grandes déviations (Proposition 7) est démontré à l'aide d'un argument de *chaîne*, celui-ci nécessite la continuité de la dérivée de la fonction  $Q$  utilisée dans le critère.

Soit  $\check{\theta}(h)$  la solution du problème de minimisation suivant :

$$(1.3.10) \quad \check{\theta}(h) = \arg \min_{t \in \Theta(M)} \tilde{m}_h(t).$$

L'estimateur de Huber local  $\check{f}^h(y)$  de  $f(y)$  est défini comme  $\check{f}^h(y) = \check{\theta}_{0,\dots,0}(h)$ .



La fonction de Huber a été introduite par [Huber \[1964,1981\]](#) dans l'idée de développer des estimateurs robustes par rapport aux valeurs extrêmes. Cette idée a été reprise en pratique par [Petrus \[1999\]](#) pour la régression gaussienne dans le cadre de la reconstruction d'images, et par [Chang et Guo \[2005\]](#) dans le modèle paramétrique gaussien pour l'estimation du positionnement par satellite. Nous proposons dans le chapitre 5 de démontrer les performances statistiques de cet estimateur. En outre, nous démontrons que la vitesse de convergence atteinte est  $n^{-\frac{\beta}{2\beta+1}}$  pour toute densité  $g_\xi$  vérifiant les hypothèses 1. [Tsybakov \[1982a\]](#) a démontré que cette vitesse est minimax pour le risque en probabilité avec la condition suivante sur  $g_\xi$ .

$$\exists p_* > 0, v_0 > 0 : \int_{\mathbb{R}} g_\xi(u) \ln \frac{g_\xi(u)}{g_\xi(u+v)} du \leq p_* v^2,$$

pour tout  $|v| \leq v_0$ . Cette condition n'est autre qu'un contrôle de la *distance de Kullback* entre les deux probabilités (Voir [Tsybakov \[2008\]](#), Chapitre 2). On peut montrer que les densités gaussienne et de cauchy vérifient cette condition. Mais par exemple, pour la régression  $\alpha$  ( $\alpha < 1/2$ ), cette condition n'est pas vérifiée.

Un cadre plus général est développé par [Arias-Castro et Donoho \[2009\]](#) pour les médianes locales itérées, justement plébiscitée pour leur capacité à traiter les sauts, et pas uniquement les zones homogènes. Leurs résultats sont aussi de type minimax pour des fonctions localement Lipschitz (dans le cadre unidimensionnel et bidimensionnel). [Arias-Castro et Donoho \[2009\]](#) montrent que le risque est du même ordre de grandeur pour les méthodes par moyennes locales ou par médianes locales dans le cadre régulier (fonctions globalement lipschitziennes). En revanche, si la fonction cible est discontinue et si les courbes régulières sont suffisamment séparées, ils montrent qu'une double itération de la méthode par médianes locales permet de diminuer la vitesse du risque par rapport au cas par moyennes ou médianes locales. Ce dernier article motive le fait que le critère  $\ell_1$  est plus performant que le critère  $\ell_2$ , notamment quand celui-ci échoue (par exemple, fonctions discontinues ou bruit sans moment).

Notons que la fonction de *Huber* utilise deux régimes, l'un reposant sur le critère  $\ell_1$  et l'autre sur le critère  $\ell_2$ . En effet, on constate que le critère  $\ell_2$  traite les petits résidus (sur l'intervalle  $[-1, 1]$ ) et le critère  $\ell_1$  les résidus plus importants notamment les valeurs extrêmes. La norme  $\ell_1$  permet d'obtenir un estimateur plus robuste (que les estimateurs linéaires par exemple). L'utilisation du critère  $\ell_2$  au voisinage de 0 implique que la dérivée de la fonction de *Huber* existe et est  $C^0$  en 0. Ce dernier point est nécessaire pour contrôler les grandes déviations (Proposition 7) de l'estimateur de Huber. Mais, on peut utiliser le critère  $\ell_2$  sur n'importe quel intervalle  $[-K, K]$  où  $K$  est un paramètre de transition entre les deux normes. En effet, il est facile d'observer que les résultats théoriques sont vérifiés pour toute constante  $K > 0$ . En pratique le problème du choix de cette constante est un problème ouvert.

L'avantage, de cette approche, est que l'on peut estimer la fonction de régression dans le modèle de *régression additive* avec densité inconnue, i.e. l'estimateur de *Huber* ne dépend pas de la densité  $g_\xi$  du bruit.

Dans le chapitre 2 (Section 2.2), nous proposons un estimateur de Huber adaptatif, construit à l'aide de la méthode de Lepski. Cette méthode nécessite un nouveau type d'inégalité de concentration de cet estimateur, développé dans le chapitre 5.

## 1.4 Mesure de l'Erreur

Après avoir développé les méthodes de construction d'estimateurs, on veut désormais choisir celui qui se sera le plus performant dans un modèle donné. Pour cela, il nous faut introduire des outils mathématiques permettant d'évaluer la performance des estimateurs et de les comparer pour choisir le meilleur. L'approche *minimax* est le premier outil développé (voir Stone [1980], Has'minskii et Ibragimov [1981] et Stone [1982]) dans ce sens pour l'estimation non-paramétrique. Pour l'adaptation, nous parlerons d'*adaptation au sens minimax* (voir Efromovich et Pinsker [1984], Lepski [1990]). Ces deux notions sont présentées dans les sections suivantes et sont utilisées tout au long de cette thèse. Notons qu'un paragraphe est consacré à l'approche *oracle* et aux *inégalités oracle*.

### 1.4.1 Approche Minimax Ponctuelle

Pour mesurer les performances des procédures d'estimation sur l'espace  $\mathbb{H}_d(\beta, L, M)$ , on utilise l'approche *minimax*. Fixons les paramètres de l'espace  $\beta, L, M > 0$ .

Rappelons que  $\mathbb{E}_f = \mathbb{E}_f^n$  est l'espérance mathématique par rapport à la loi de probabilité  $\mathbb{P}_f = \mathbb{P}_f^n$  des observations  $\mathcal{Z}_n$  satisfaisant le modèle de régression *générale* (1.2.2). Premièrement nous définissons le risque maximal sur  $\mathbb{H}_d(\beta, L, M)$  correspondant à l'estimation de la fonction  $f$  en un point donné  $y \in [0, 1]^d$ . On parlera d'*estimation ponctuelle*. Deuxièmement, on définit le risque minimax, et enfin les notions, de vitesse de convergence minimax et estimateurs minimax, sont données.

Soit  $\tilde{f}$  un estimateur construit à partir des observations  $\mathcal{Z}_n$  et choisi de façon arbitraire.

**Définition 6.** Soit  $q > 0$ , on appelle *risque maximal de l'estimateur  $\tilde{f}$  sur  $\mathbb{H}_d(\beta, L, M)$*  la quantité suivante

$$R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M)] = \sup_{f \in \mathbb{H}_d(\beta, L, M)} \mathbb{E}_f |\tilde{f}(y) - f(y)|^q.$$

En général, on s'intéresse à la vitesse de convergence  $\psi_n$  vers 0 du risque maximal d'un estimateur  $\tilde{f}$ , nous cherchons un estimateur qui atteint la "meilleure" vitesse de convergence pour le risque maximal. Remarquons que cette vitesse sera donnée pour la pire fonction à estimer dans la classe de Hölder  $\mathbb{H}_d(\beta, L, M)$ , cette approche est souvent considérée comme pessimiste, mais elle est nécessaire lorsque l'on suppose que la fonction  $f$  appartient à un espace fonctionnel. Il est bon de noter que le risque maximal est utilisé dans cette thèse avec la semi-norme ponctuelle. L'avantage principale de ce risque est de choisir de façon locale la fenêtre d'estimation. Donc on aura un estimateur qui s'adapte localement à la régularité de

la fonction  $f$  à estimer. Les résultats de type borne supérieure pour ce risque sont souvent les plus techniques à établir, il est possible d'obtenir le risque “global” en norme  $L_p$  en intégrant le risque “ponctuel” et en utilisant le théorème de Fubini.

**Définition 7.** On appelle *risque minimax* sur  $\mathbb{H}_d(\beta, L, M)$  la quantité suivante

$$R_{n,q}[\mathbb{H}_d(\beta, L, M)] = \inf_{\hat{f}} R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M)],$$

où  $\inf$  est pris sur l'ensemble des estimateurs (fonctions mesurables des observations).

Cette quantité peut être décrite de la façon suivante : “on mesure de manière théorique le risque du meilleur estimateur possible pour la pire fonction à estimer dans une boule de Hölder”.

**Définition 8.** La suite de normalisations  $\psi_n(\beta)$  est appelée *vitesse de convergence minimax* et l'estimateur  $\hat{f}$  est dit *minimax* (asymptotiquement minimax) si

$$\begin{aligned} \liminf_{n \rightarrow \infty} \psi_n^{-q}(\beta) R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M)] &> 0; \\ \limsup_{n \rightarrow \infty} \psi_n^{-q}(\beta) R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M)] &< \infty. \end{aligned}$$

Pour un modèle donné, notre but est de construire l'estimateur minimax (celui qui atteint la vitesse minimax). Techniquement, il nous faut montrer une borne supérieure du risque maximal et une borne inférieure pour le risque minimax. Ces deux bornes doivent être égales à constante près. Pour la borne inférieure, de nombreuses techniques ont été mises en place, on renvoie le lecteur au livre de [Tsybakov \[2008\]](#) pour une présentation de ces méthodes. Plus récemment, dans l'article [Chichignoud \[2010a\]](#) (Chapitre 4, Section 4.2), nous utilisons le critère introduit par [Klutchnikoff \[2005\]](#) pour obtenir la borne inférieure du risque minimax (voir Chapitre 4). Pour les bornes supérieures, la forme de l'estimateur en question et des techniques classiques sont utilisées comme la décomposition biais-variance.

Dans la majorité des cas, les estimateurs minimax dépendent des paramètres de régularité  $\beta$  et  $L$ , ( $\hat{f} = \hat{f}_{\beta,L}$ ). En pratique, le problème est que les paramètres de régularité  $\beta$ ,  $L$  ne sont pas connus par le statisticien. De plus, ces paramètres ne peuvent pas être estimés au préalable à l'inverse de la borne supérieure  $M$ . Dans ce sens, l'*adaptation* a été développée dans les années 90 par [Efromovich et Pinsker \[1984\]](#), [Lepski \[1990\]](#), [Donoho et Johnstone \[1994\]](#), [Barron, Birgé, et Massart \[1999\]](#), etc. Ces méthodes permettent de construire des estimateurs qui ne dépendent pas de la régularité de la fonction cible et qui ont une “bonne” performance.

**Approche Maxiset.** Plus récemment, la notion de *maxiset* a été introduite par [Kerkycharian et Picard \[2002\]](#) et largement développée par [Autin \[2006\]](#), [Autin, Picard, et Rivoirard \[2006\]](#) et [Autin \[2008\]](#). Cette approche est complémentaire aux motivations énoncées

dans ce paragraphe. Etant donnés un estimateur  $\hat{f}$  et une vitesse de convergence  $\varphi_n$ , l'idée est de chercher l'espace fonctionnel le plus grand possible sur lequel  $\hat{f}$  atteint la vitesse  $\varphi_n$ . Cette notion est beaucoup moins pessimiste que le risque minimax puisqu'un estimateur n'est plus défini par "la plus mauvaise" performance sur un espace préalablement fixé, mais par l'ensemble des fonctions pouvant être estimées par  $\hat{f}$  à la vitesse  $\varphi_n$ . Cette approche peut être utilisée pour une large gamme de modèles. Cette approche, qui peut être intéressante pour les praticiens, dans le sens où elle exhibe les fonctions "bien estimées" (à vitesse choisie) par un estimateur choisi. L'inconvénient est que deux estimateurs ne sont pas toujours comparables si leurs maxisets ne sont pas emboîtés (par exemple, seuillage par blocs et seuillage par arbre).

### 1.4.2 Approche Minimax Adaptative

Dans cette section, nous présentons les contraintes de l'adaptation et les résultats attendus. Le but de l'adaptation au sens minimax est de choisir un estimateur  $\hat{f}^*$  qui ne dépend pas de  $\beta, L$  et qui atteint la "meilleure" vitesse de convergence pour tout  $\beta, L$ . Nous verrons qu'il est souvent nécessaire de payer un prix pour l'adaptation, c'est-à-dire que l'on peut construire un estimateur qui ne dépend pas de  $\beta, L$  mais qui atteint seulement la vitesse  $\psi_n^p$  minimax pénalisée. La question que l'on se pose naturellement est la suivante : Existe-t-il un estimateur  $\hat{f}^*$  qui ne dépend pas de  $\beta, L$  et qui atteint la vitesse minimax  $\psi_n$  ? Pour les risques  $L_p$  la réponse à cette question est positive, voir [Efremovich et Pinsker \[1984\]](#), [Lepski \[1991\]](#), [Donoho, Johnstone, Kerkycharian, et Picard \[1995\]](#), [Lepski, Mammen, et Spokoiny \[1997\]](#), [Goldenshluger et Nemirovski \[1997\]](#), [Juditsky \[1997\]](#) et [Goldenshluger et Lepski \[2008\]](#). Pour le risque ponctuel, le prix à payer n'est pas nul, [Lepski \[1990,1992a,1992b\]](#), [Brown et Low \[1996\]](#), [Spokoiny \[1996\]](#), [Lepski et Spokoiny \[1997\]](#), [Tsybakov \[1998\]](#), [Klutchnikoff \[2005\]](#), [Chichignoud \[2010a,2010b,2010c\]](#). En général, le prix à payer est une puissance de  $\ln n$ .

**Définition 9.** On appelle  $\{\psi_n^p(\beta)\}_{\beta,L}$  une famille de normalisations indexée par les paramètres de régularité. On dit que l'estimateur  $\hat{f}^*$  est  $\psi$ -adaptatif si et seulement si  $\forall \beta, L$  on a

$$\limsup_{n \rightarrow \infty} (\psi_n^p(\beta))^{-q} R_{n,q}[\hat{f}^*, \mathbb{H}_d(\beta, L, M)] < \infty.$$

Remarquons que la vitesse de convergence ne dépend pas de  $L$ . C'est en effet souvent le cas : la constante de Lipschitz intervient de façon polynômiale dans les constantes de la borne supérieure.

La recherche de la meilleure famille de normalisations (celle qui tend le plus vite vers 0) est une priorité. L'optimalité de ces familles a été étudiée par [Lepski \[1992a\]](#), [Lepski \[1992b\]](#), [Lepski et Spokoiny \[1997\]](#), [Tsybakov \[1998\]](#), [Klutchnikoff \[2005\]](#) et [Chichignoud \[2010a\]](#). Pour obtenir les bornes supérieures, le recours à des inégalités de concentration exponentielles est nécessaire (voir Section 1.5). On renvoie le lecteur au livre de [Massart \[2007\]](#) sur les inégalités maximales et de concentration. Rappelons que notre estimateur  $\hat{f}^*$  est *adaptatif* pour la pire fonction dans la boule de Hölder. Ce point de vue est souvent considéré comme pessimiste.

**Approche Oracle.** Il est possible de construire des méthodes adaptatives à toutes fonctions dans un espace de Hilbert  $H$ . Etant données une méthode  $\hat{f}$  et une famille d'estimateurs  $\Lambda$ , une inégalité du type :

$$\mathbb{E}_f \|\hat{f} - f\|^2 \leq (1 + \rho_n) \inf_{\tilde{f} \in \Lambda} \mathbb{E}_f \|\tilde{f} - f\|^2 + T_n, \quad \forall f \in H,$$

et  $\rho_n \geq 0$  est appelée **inégalité oracle**.

On appelle *oracle* l'estimateur  $\hat{f}_0 = \arg \min_{\tilde{f} \in \Lambda} \mathbb{E}_f \|\tilde{f} - f\|^2$ . Connaissant la fonction  $f$ ,  $\hat{f}_0$  correspond au meilleur estimateur possible dans la famille  $\Lambda$ .  $T_n$  est un terme résiduel à contrôler.

Dans de nombreuses situations, une inégalité oracle est un très bon moyen d'obtenir des résultats d'adaptation au sens minimax, i.e. en utilisant le fait que l'oracle dans la famille est un estimateur minimax (résultat à démontrer au préalable), on peut montrer que l'estimateur  $\hat{f}$  est adaptatif optimal (atteint la vitesse minimax et ne dépend pas de la régularité de l'espace considéré). Noter qu'en général ceci est possible seulement dans le cas de l'estimation globale (en norme  $L_p$ ), les inégalités d'oracle pour la semi-norme ponctuelle n'apparaissent pas dans la littérature, notamment pour leur différence (prix à payer pour l'adaptation) avec la norme  $L_p$ . Ce type de résultat est cependant beaucoup plus fort, il est établi sans aucune restriction sur la fonction  $f$ . Par ailleurs, une inégalité oracle présente un caractère non-asymptotique dans la mesure où elle permet de comparer les performances d'une méthode à un estimateur "idéal"  $\hat{f}_0$ , et ce, quelque soit le nombre d'observations.

L'estimation en grande dimension est aussi très étudiée aujourd'hui. En outre, la notion de *parcimonie* (en anglais : *sparsity*) est systématiquement utilisée pour obtenir des résultats du type *inégalités oracle* (voir par exemple, [Bunea, Tsybakov, et Wegkamp \[2007\]](#)).

Dans cette thèse, les résultats présentés sont asymptotiques par souci de présentation. Néanmoins, au vu des preuves, il est possible d'obtenir des résultats non-asymptotiques, c'est-à-dire que notre borne supérieure sera composée de la vitesse  $\psi_n^p$ , plus un reste  $T_n$ . En général, on choisit  $n$  assez grand pour que le reste soit négligeable. Nous verrons dans la section 1.5.2 que la méthode de Lepski n'a nullement un caractère asymptotique.

Pour l'approche *oracle*, rappelons que la norme  $\|\cdot\|$  est une norme  $L_p$ . Il est souvent difficile d'obtenir des inégalités oracle avec une semi-norme ponctuelle. Seuls les travaux de [Plancade \[2008\]](#), sur l'estimation de la densité du bruit dans le modèle de régression homoscédastique, vont dans ce sens. Nos résultats sont tous donnés avec le risque ponctuel. Nous verrons que la principale difficulté est d'obtenir un contrôle des grandes déviations avec la semi-norme ponctuelle.

## 1.5 Adaptation

Cette partie contient un panorama sur les méthodes adaptatives développées ces vingt dernières années. Nous utilisons, dans cette thèse, la méthode de Lepski pour l'adaptation.

Une section y est consacrée pour présenter l'idée générale et les conditions nécessaires à l'application de la méthode.

### 1.5.1 Généralités

Une des premières méthodes *adaptatives* a été développée par Lepski [1990]. Nous la présentons dans la section 1.5.2. On citera, sans plus de détails, une autre approche possible pour l'adaptation, les méthodes dites *d'agrégation d'estimateurs* (voir notamment Nemirovski [2000], Rigollet et Tsybakov [2007] et Lecué et Mendelson [2009]). Dans la suite, nous présentons deux méthodes d'adaptation très populaires : le *seuillage par ondelettes* et la *sélection de modèles*.

**Seuillage par ondelettes.** Suite à la création des ondelettes par Meyer [1992], les premières méthodes adaptatives fondées sur les ondelettes sont apparues sous l'impulsion de Donoho et Johnstone [1994] et Donoho, Johnstone, Kerkycharian, et Picard [1995]. L'idée est de projeter la fonction de régression dans une base d'ondelettes et d'estimer les coefficients. On procède par un seuillage (déterministe ou aléatoire) des coefficients pour ne retenir que les plus grands coefficients à savoir ceux qui contiennent l'information. D'un point de vue maxiset, les méthodes de seuillage par arbre ou par blocs ont été jugées très performantes par Autin [2004]. En particulier, leur maxiset contient les espaces de Besov.

**Sélection de modèles.** On souhaite construire une procédure établie à partir des observations  $(Y_i)_i$  qui permet de choisir un estimateur parmi la collection  $\{\hat{f}_m : m \in \mathcal{M}\}$  où  $\hat{f}_m$  est un minimiseur de  $\gamma_n$  sur  $S_m$ ,  $\{S_m : m \in \mathcal{M}\}$  une collection de modèles. L'objectif est de choisir un estimateur dont le risque est aussi proche que possible de celui de l'oracle. Pour faire ce choix, on procède par pénalisation. On choisit l'estimateur adaptatif  $\hat{f}_{\hat{m}}$  tel que

$$\hat{m} = \arg \min \left[ \gamma_n(\hat{f}_m) + \text{pen}(m) \right],$$

où  $\gamma_n(\cdot)$  est un contraste empirique et  $\text{pen}(m)$  est une pénalité à choisir. Les premiers résultats obtenus par critère pénalisé sont dûs à Akaike [1973] pour l'estimation de densité par vraisemblance pénalisée et à Mallows [1973] pour l'estimation de la fonction de régression dans un cadre gaussien. Pour la pénalité  $\text{pen}_{\text{Akaike}}(m) = D_m$  où  $D_m$  est la dimension de  $S_m$ , Akaike a obtenu le critère AIC et, pour la pénalité  $\text{pen}_{\text{Mallows}}(m) = 2\sigma D_m$ , le critère s'appelle  $C_p$  de Mallows. Birgé et Massart [2001] ont montré que ces critères sont asymptotiquement optimaux à condition que la taille de  $\mathcal{M}$  ne soit pas trop grande. Dans un cadre non-asymptotique, Birgé et Massart [2001] ont introduit des pénalités plus générales que celle de Mallows

$$\text{pen}_{\text{bm}} = C \sigma^2 D_m \left( 1 + \sqrt{2L_m} \right)^2, \quad m \in \mathcal{M},$$

où  $C > 1$  et  $\{L_m : m \in \mathcal{M}\}$  une collection de poids positifs. Birgé et Massart ont obtenu un contrôle sur le risque de  $f$  quelque soit la taille de la collection de modèles. La construction de leur procédure se base sur des résultats de concentration de la mesure gaussienne. La qualité d'une procédure de sélection de modèle correspond à sa capacité à choisir un estimateur  $\tilde{f}$  parmi une collection  $\{\hat{f}_m : m \in \mathcal{M}\}$  qui ait un risque faible. Afin d'évaluer cette qualité, on peut considérer les inégalités de type oracle pour valider théoriquement la qualité de l'estimateur  $\tilde{f}$ . Il nous faut donc établir des inégalités de la forme

$$\mathbb{E}_f l(f, \tilde{f}) \leq K \inf_{m \in \mathcal{M}} \left\{ \inf_{f' \in S_m} l(f, f') + \text{pen}(m) \right\} + R, \quad \forall f,$$

où  $l(., .)$  est une fonction de perte,  $K$  est une constante universelle et  $R$  un terme de reste. En général, la fonction de perte est celle des moindres carrées ou la divergence de Kullback-Leibler. Le principal avantage est qu'elle n'impose aucune restriction sur la fonction à estimer  $f$ . En pratique, cette méthode pose la question du choix de la constante  $C$ , qui est une constante de régularité.

Plus récemment, les travaux de [Cavalier et Golubev \[2006\]](#) et [Marteau \[2009\]](#) utilisent une méthode dite *d'enveloppe du risque*, fondée sur la sélection de modèles, pour les problèmes inverses très mal-posés.

### 1.5.2 Choix de la Fenêtre : Méthode de Lepski

Comme nous l'avons vu, les méthodes utilisant une fenêtre fixe ne peuvent rendre compte de la géométrie (ou de régularité) de la fonction. Une amélioration possible est alors de choisir le voisinage de manière locale, en s'appuyant sur les données observées. L'idée de ce type d'approche est de sélectionner le voisinage le plus pertinent possible à partir des observations. Cette idée bien connue sous le nom de méthode dite de *Lepski* ([Lepski \[1990\]](#)) permet de sélectionner la fenêtre de façon adaptative parmi une grille, cette méthode repose sur la comparaison d'estimateurs appartenant à une grille finie. Cette opération est souvent très chère en complexité de calcul.

**Chronologie.** La première procédure adaptative fût donnée par [Efromovich et Pinsker \[1984\]](#). [Stone \[1982\]](#) avait posé, sans la résoudre, la question de l'existence d'un estimateur ne dépendant pas de la régularité de la fonction cible. [Lepski \[1990\]](#) donna une nouvelle procédure qui portera plus tard son nom. Il construit un estimateur adaptatif pour choisir la régularité de Hölder dans le modèle du bruit blanc gaussien pour l'estimation ponctuelle. Il démontra qu'il n'existe pas d'estimateur adaptatif efficace. En effet, il est fréquent de payer un prix  $\sqrt{\ln \varepsilon^{-1}}$  pour l'adaptation. C'est-à-dire que l'on ne peut pas faire mieux que la vitesse  $(\varepsilon \sqrt{\ln \varepsilon^{-1}})^{\frac{2\beta}{2\beta+1}}$ .

Dans l'article [Lepski \[1991\]](#), une procédure plus générale est donnée pour n'importe quel type d'estimateurs d'une fonction. On peut trouver dans cet article les conditions suffisantes



pour obtenir un estimateur adaptatif efficace pour le risque  $L_p$ . Dans les deux articles [Lepski \[1992a, 1992b\]](#), on peut trouver des exemples où il n'existe pas d'estimateurs efficaces. Cela dépend généralement de la fonction de perte utilisée dans le risque (norme  $L_p$ , semi-norme ponctuelle ou encore norme infinie). Ces trois articles traitent entre autres le modèle du bruit blanc, d'estimateurs à noyau et d'estimateurs par polynômes locaux.

Plusieurs articles ont ensuite utilisé la méthode de Lepski pour l'adaptation et ont contribué à l'amélioration de celle-ci. Un choix de la fenêtre est proposé dans [Lepski, Mammen, et Spokoiny \[1997\]](#) pour les estimateurs à noyau. Ils obtiennent sur les espaces de Besov les vitesses adaptatives efficaces, qui étaient déjà connues par [Donoho, Johnstone, Kerkycharian, et Picard \[1995\]](#) avec les ondelettes. À noter que dans cet article, le choix de la fenêtre est fait avec une méthode ponctuelle et les vitesses sont obtenues en norme  $L_p$ .

La régression gaussienne à design fixe a fait l'objet d'un article de [Goldenshluger et Nemirovski \[1997\]](#), où l'estimateur des polynômes locaux est utilisé avec le risque ponctuel et  $L_p$  sur les espaces de Sobolev. [Spokoiny \[1998\]](#) a étudié les fonctions discontinues et propose une méthode adaptative reposant sur l'idée de Lepski.

[Lepski et Levit \[1999\]](#) et [Kerkycharian, Lepski, et Picard \[2001\]](#) proposent des méthodes adaptatives pour les fonctions anisotropes. La nouveauté est dans le caractère multidimensionnel de la fenêtre d'estimation, avec l'introduction d'un estimateur artificiel dans la procédure de Lepski. On note la présence d'inégalités oracle dans [Kerkycharian, Lepski, et Picard \[2001\]](#) avec le risque  $L_p$ . Citons les travaux de [Klutchnikoff \[2005\]](#) sur l'estimation ponctuelle de fonctions de Hölder anisotropes. Plus récemment, l'article de [Goldenshluger et Lepski \[2008\]](#) donne une procédure s'appuyant sur l'idée de Lepski pour un large choix d'estimateurs linéaires sans condition sur le modèle (régression gaussienne, bruit blanc gaussien et modèle de densité). [Goldenshluger et Lepski \[2009a\]](#) donnent une procédure qui permet d'obtenir des inégalités oracle sur une très large famille d'estimateurs linéaires.

Il est possible d'utiliser la règle de Lepski pour les méthodes de seuillage des coefficients d'ondelettes. En effet, [Autin \[2004\]](#) a mis en place une procédure dite *procédure hard tree*, qui utilise la règle de Lepski pour choisir les coefficients à seuiller dans l'arbre. Dans le modèle du bruit blanc gaussien, cette méthode permet d'avoir un maxiset plus important que pour les méthodes de seuillage classiques.

La méthode de Lepski a aussi été utilisée pour les problèmes inverses. L'article de [Mathé \[2006\]](#) est une très bonne description de la méthode dans le cas des problèmes inverses déterministes et aléatoires. Plus récemment, on peut citer une nouvelle procédure de test introduite par [Spokoiny et Vial \[2009\]](#) pour les problèmes inverses.

En pratique, on retrouve souvent la règle de Lepski du fait de son implémentation assez simple, malgré le coût important de la complexité. En effet, il est nécessaire de calculer tous les estimateurs dans la famille que nous considérons. Les travaux de [Polzehl et Spokoiny \[2000, 2003\]](#) ont considéré la méthode de Lepski pour le débruitage d'images. L'estimateur du maximum de vraisemblance local est utilisé pour l'adaptation dans les modèles dits *réguliers* dans les articles de [Polzehl et Spokoiny \[2006\]](#) et [Katkovnik et Spokoiny \[2008\]](#), notamment pour les modèles poissonniens adaptés aux images. Dans le même temps, [Katkovnik \[1999\]](#) et,



par la suite [Astola, Egiazarian, et Katkovnik \[2002\]](#) et [Astola, Egiazarian, Foi, et Katkovnik \[2010\]](#), donnent une très bonne illustration de comment utiliser la règle de Lepski en pratique. Par exemple, ils introduisent une nouvelle interprétation de cette méthode que [Katkovnik \[1999\]](#) nomme *ICI* (en anglais : *Intersection of Confidence Intervals*) et utilisent des voisinages anisotropes circulaires, ce qui améliore le débruitage de l'image. Cette méthode, explicitée plus loin, permet de réduire considérablement la complexité de la méthode tout en gardant ces bonnes propriétés théoriques.



FIGURE 1.9 – Illustration de l'effet poivre et sel sur une image traitée avec la méthode de [Kervrann et Boulanger \[2006\]](#) (voir par exemple le coin en bas à droite)

[Kervrann et Boulanger \[2006\]](#) ont utilisé la méthode de Lepski pour le choix de la zone de recherche pour les *Moyennes Non Locales* (en anglais : Non-Local Means ou encore NL-Means, voir [Buades, Coll, et Morel \[2005\]](#)). Ces articles mettent en avant un inconvénient, en pratique, de la méthode : l'image obtenue après traitement a une tendance à être dégradée par un bruit de type poivre et sel (Figure 1.9). Ceci est dû à la procédure, qui avec une probabilité faible, s'arrête trop tôt (voir Description). Pour combler cette défaillance, l'utilisation de filtres médians est souvent préférée à celle des filtres linéaires (voir [Astola, Egiazarian, Foi, et Katkovnik \[2010\]](#)).

L'idée de cette section est de comprendre comment on peut utiliser la méthode de Lepski, donner les idées de la preuve, étayées par quelques calculs, et expliquer les conditions

suffisantes à rassembler pour garantir la performance de la procédure.

**Description.** Comme nous l'avons précisé ci-dessus, l'approche développée par Lepski [1991] est conçue pour l'adaptation au sens minimax. On peut cependant développer une approche oracle (voir Goldenshluger et Lepski [2009a]). Dans cette partie, nous considérons la régression généralisée définie dans la section 1.2.2 et l'estimateur localement paramétrique  $\hat{f}^h(y) = \hat{\theta}_{0,\dots,0}(h)$  tel que

$$\hat{\theta}(h) = \arg \min_{t \in \Theta(M)} \delta(t, \mathcal{Z}_n),$$

où  $\delta(.,.)$  est un contraste empirique. Par exemple,  $\delta(.,.)$  peut être le critère des moindres carrés ou bien le critère bayésien (1.3.7) ou encore le critère de Huber (1.3.9). Rappelons que  $f$  appartient à une boule de Hölder  $H_d(\beta, L)$  (voir Définition 1), où  $\beta$  et  $L$  sont inconnus pour le statisticien.

Supposons que nous connaissons déjà la vitesse minimax que nous noterons  $\psi_{n,\gamma}(\beta) = n^{-\frac{\beta}{\gamma\beta+d}}$ ,  $1 \leq \gamma \leq 2$ . Nous verrons que, dans le modèle de régression gaussienne, la vitesse est  $\psi_{n,2}(\beta)$  et dans le modèle de régression multiplicative uniforme la vitesse est  $\psi_{n,1}(\beta)$  (voir Section 3.4). Nous supposons aussi connue la construction de l'estimateur minimax noté  $\hat{f}_{h(\beta,L)}$  où  $h(\beta, L)$  est l'argument qui permet d'équilibrer le biais et la variance de notre estimateur que l'on peut écrire comme la solution du problème de minimisation suivant :

$$(1.5.1) \quad h(\beta, L) = (L^\gamma n)^{-\frac{1}{\gamma\beta+d}} = \arg \min_{h>0} Ldh^\beta + (nh^d)^{-\frac{1}{\gamma}}.$$

L'estimateur  $\hat{f}_{h(\beta,L)}$  dépend explicitement de  $\beta, L$  à travers la fenêtre  $h(\beta, L)$ . Choisir la régularité de façon adaptative revient alors à choisir la fenêtre. La procédure que nous présentons est celle développée par Lepski, Mammen, et Spokoiny [1997] que l'on appelle abusivement *méthode de Lepski*. Cette procédure a été donnée spécialement pour le choix de la fenêtre, et donc pour les estimateurs localement paramétriques. Néanmoins, l'idée peut être utilisée dans un autre contexte. Spokoiny et Vial [2009] utilisent le principe de Lepski pour les problèmes inverses. Leur estimateur est la somme successive des coefficients bruitées de la fonction cible. Cet estimateur s'arrête au  $N$ -ième coefficient. La question est de savoir à quel rang s'arrêter. Les auteurs développent une méthode pour tester si l'on peut passer au rang supérieur ou non. Ceci est fait à partir de l'idée de Lepski [1990].

L'idée de l'adaptation est la suivante : on dispose d'une famille d'estimateurs  $\{\hat{f}^h\}_{h \in \mathcal{H}_n}$ , où  $\mathcal{H}_n \subseteq [0, 1]$  tel que  $h(\beta, L) \in \mathcal{H}_n$ . Cette famille contient un oracle, ici l'estimateur minimax  $\hat{f}_{h(\beta,L)}$ . Nous supposons inconnues les quantités  $\beta$  et  $L$ , et on admet connaître l'information supplémentaire que  $\beta \in ]0, b]$  et  $L > 0$  où  $b \in \mathbb{N}^*$  est connu. Pour la constante  $L$ , aucune restriction n'est imposée, en effet, son rôle est minimisé dans le choix de la fenêtre, du fait qu'elle n'intervient pas dans la vitesse de convergence.

Notons  $\mathcal{H}_n = [h_{\min}, h_{\max}]$  où  $h_{\min} < h_{\max}$  sont à choisir pour que  $h(\beta, L) \in \mathcal{H}_n$ . La procédure de Lepski commence par la construction d'une grille sur l'ensemble  $\mathcal{H}_n$  des

fenêtres, soit

$$(1.5.2) \quad h_k = 2^{-k} h_{\max}, \quad k = 0, \dots, \mathbf{k}_n,$$

où  $\mathbf{k}_n$  est le plus grand entier tel que  $h_{\mathbf{k}_n} \in \mathcal{H}_n$ . Soit

$$(1.5.3) \quad \mathcal{F} = \left\{ \mathring{f}^{(k)}(y) = \mathring{f}^{h_k}(y), \quad k = 0, \dots, \mathbf{k}_n \right\},$$

la grille construit sur notre famille d'estimateurs localement paramétriques.

**Objectif.** Le but de l'adaptation est de choisir un estimateur  $\mathring{f}^{(\hat{k})} \in \mathcal{F}$  où  $\hat{k}$  est une fonction mesurable des observations  $\mathcal{Z}_n$ . Une fois la procédure établie, il nous faut vérifier que notre estimateur est bien *adaptatif* (voir Définition 9). C'est-à-dire qu'il atteint simultanément sur chaque espace  $\mathbb{H}_d(\beta, L)$  la vitesse  $\psi_n^p(\beta)$ ,  $\forall \beta \in [0, b[$ . Notons que cette famille de normalisations doit être optimale (voir par exemple Chapitre 4). On peut écrire la vitesse adaptative de la façon suivante :

$$\psi_n^p(\beta) = \left( n^{-1} \rho_n(\beta) \right)^{\frac{\beta}{\gamma\beta+d}},$$

où  $\rho_n(\beta)$  est un terme dit de *pénalité* ou de *prix à payer* pour l'adaptation. En général, pour l'estimation ponctuelle, on choisit la pénalité égale à  $[1 + (b - \beta) \ln n]^{1/\gamma}$  (voir Lepski [1990,1992a,1992b], Brown et Low [1996], Spokoiny [1996], Lepski et Spokoiny [1997], Tsybakov [1998], Klutchnikoff [2005], Chichignoud [2010a,2010b,2010c]), la description de la procédure permettra d'expliquer ce choix. En norme  $L_p$ , la pénalité est souvent égale à 1 (voir Efromovich et Pinsker [1984], Lepski [1991], Donoho, Johnstone, Kerkycharian, et Picard [1995], Lepski, Mammen, et Spokoiny [1997], Goldenshluger et Nemirovski [1997], Juditsky [1997] et Goldenshluger et Lepski [2008]).

### Règle de sélection pour deux fenêtres.

Pour commencer considérons la méthode dans le cas simple du choix entre deux estimateurs. Notons les  $\mathring{f}^{(1)} = \mathring{f}^{h_1}$ ,  $0 < h_1 < 1$  et  $\mathring{f}^{(2)} = \mathring{f}^{h_2}$ ,  $h_1 < h_2 < 1$ , on est dans le cas où  $\text{Card}(\mathcal{F}) = 2$ . On dispose de l'information supplémentaire : “un des deux estimateurs est minimax sur  $\mathbb{H}_d(\beta, L)$ ” par construction de  $\mathcal{F}$ . Nous allons construire une procédure pour choisir un estimateur parmi les deux pour atteindre la vitesse adaptative. Soit la règle de sélection suivante :

$$f^* = \begin{cases} \mathring{f}^{(2)}, & |\mathring{f}^{(1)}(y) - \mathring{f}^{(2)}(y)| \leq C \sigma_p(h_1) \\ \mathring{f}^{(1)}, & |\mathring{f}^{(1)}(y) - \mathring{f}^{(2)}(y)| > C \sigma_p(h_1) \end{cases}$$

Nous appelons *écart-type*  $\sigma_h = (nh^d)^{-1/\gamma}$  et *écart-type pénalisé*  $\sigma_p(h) = \text{pen}_h \times \sigma(h)$  avec  $\text{pen}_h$  une pénalisation à choisir en fonction du problème.  $C$  est une constante que l'on peut voir comme un paramètre de seuillage de la méthode. Nous choisirons celle-ci plus tard. Nous pouvons garantir la performance de notre estimateur par le résultat suivant.

**Proposition 1.** *Pour tout  $\beta, L \in [0, b[ \times \mathbb{R}_+$  on a*

$$\sup_{f \in \mathbb{H}_d(\beta, L)} \mathbb{E}_f |f^*(y) - f(y)| \leq (2c_1 + c_2) \psi_n^p(\beta) + R_n,$$

où  $c_1, c_2$  sont des constantes positives et  $R_n$  est un terme de reste que l'on définira dans la suite.

Avec ce proposition, l'estimateur  $f^*$  atteint la vitesse  $\psi_n^p(\beta)$  pour tout  $\beta, L \in [0, b[ \times \mathbb{R}_+$ .

**Idée de la preuve.** Soit  $h^p(\beta, L) = \left( \frac{\rho_n(\beta)}{L^{\gamma_n}} \right)^{-\frac{1}{\gamma_{\beta+d}}}$  la fenêtre minimax pénalisée. Dans ce paragraphe, nous expliquons comment fonctionne la méthode de Lepski. On peut comprendre comment choisir la vitesse pénalisée, la constante  $C$  et la variance pénalisée.

1er cas :  $\hat{f}^{(1)}$  est minimax pénalisé, c'est-à-dire que  $h_1 = h^p(\beta, L)$ . Définissons l'événement

$$A = \left\{ |\hat{f}^{(1)}(y) - \hat{f}^{(2)}(y)| \leq C \sigma_p(h_1) \right\}.$$

**1)** Si  $A^c$  est réalisé alors la règle de sélection donne  $f^* = \hat{f}^{(1)}$  (la règle a choisi le bon estimateur). Comme  $\hat{f}^{(1)}$  est minimax, on a

$$\mathbb{E}_f |f^*(y) - f(y)| \mathbb{I}_{A^c} \leq \mathbb{E}_f |\hat{f}^{(1)}(y) - f(y)| \leq c_1 \psi_n^p(\beta).$$

**2)** Si  $A$  est réalisé alors  $f^* = \hat{f}^{(2)}$  (ici la procédure s'est trompée d'estimateur), nous allons regarder quelle est la perte due à cette erreur.

$$\mathbb{E}_f |f^*(y) - f(y)| \mathbb{I}_A = \mathbb{E}_f |\hat{f}^{(2)}(y) - f(y)| \mathbb{I}_A \leq \mathbb{E}_f |\hat{f}^{(2)}(y) - \hat{f}^{(1)}(y)| \mathbb{I}_A + \mathbb{E}_f |\hat{f}^{(1)}(y) - f(y)| \mathbb{I}_A.$$

Sous l'événement  $A$ , le premier terme est contrôlé par  $C \sigma_p(h_1)$  et comme  $\hat{f}^{(1)}$  est minimax, il existe une constante  $c_1$  telle que

$$\mathbb{E}_f |f^*(y) - f(y)| \mathbb{I}_A \leq C \sigma_p(h_1) + c_1 \psi_n^p(\beta).$$

Nous pouvons montrer que si  $\text{pen}_{h^p(\beta, L)} = \rho_n(\beta)$  alors il existe une constante  $c_2$  telle que

$$\sigma_p(h_1) = \sigma_p(h^p(\beta, L)) \leq c_2 \psi_n^p(\beta),$$

ainsi on obtient avec les deux inégalités ci-dessus.

$$\mathbb{E}_f |f^*(y) - f(y)| \mathbb{I}_A \leq (c_1 + c_2) \psi_n^p(\beta).$$

2ème cas :  $\hat{f}^{(2)}$  est minimax pénalisé ( $h_2 = h^p(\beta, L)$ ). Dans ce cas, la probabilité que l'événement  $A^c$  se réalise est faible.

1) Si  $A$  est réalisé alors  $f^* = \mathring{f}^{(2)}$  (la procédure donne le bon estimateur). Ainsi

$$\mathbb{E}_f |f^*(y) - f(y)| \mathbb{I}_A \leq \mathbb{E}_f |\mathring{f}^{(2)}(y) - f(y)| \leq c_1 \psi_n^p(\beta).$$

2) Si  $A^c$  est réalisé alors la règle de sélection donne  $f^* = \mathring{f}^{(1)}$  (la procédure s'est trompée). En utilisant l'inégalité de Cauchy-Schwarz, on obtient :

$$\mathbb{E}_f |f^*(y) - f(y)| \mathbb{I}_{A^c} \leq \left( \mathbb{E}_f |\mathring{f}^{(1)}(y) - f(y)|^2 \right)^{1/2} \sqrt{\mathbb{P}_f \{A^c\}}.$$

Le premier facteur (risque de l'estimateur  $\mathring{f}^{(1)}$ ) tend vers 0 (car c'est un estimateur de la fonction  $f$  avec une mauvaise fenêtre) mais pas à la bonne vitesse. Il est facile de montrer que ce facteur est néanmoins borné par une constante  $K$ . Dans certains exemples, on peut le calculer explicitement et connaître sa vitesse. Le second terme  $\sqrt{\mathbb{P}_f \{A^c\}}$  va compenser la perte. Nous pouvons donner une condition suffisante pour que la règle de sélection à deux fenêtres fonctionne. S'il existe une constante  $c_3$ , une pénalité  $\text{pen}_h$ , une constante  $C$  et une contante  $a = a(C) > 1$  telles que  $\text{pen}_{h^p(\beta, L)} = \rho_n(\beta)$  et

$$(1.5.4) \quad \left( \mathbb{P}_f \left\{ |\mathring{f}^{(1)}(y) - \mathring{f}^{(2)}(y)| > C \sigma_p(h_1) \right\} \right)^{1/2} \leq c_3 n^{-a}.$$

Alors d'après les cas traités ci-dessus, on a

$$\mathbb{E}_f |f^*(y) - f(y)| \leq (2c_1 + c_2) \psi_n^p(\beta) + K c_3 n^{-a}, \quad \forall \beta, L \in [0, b[ \times \mathbb{R}_+.$$

où  $R_n = K c_3 n^{-a}$  est le terme de reste. Le proposition 1 est démontré.

Nous pouvons donner un schéma qui permet de contrôler la probabilité (1.5.4). Rappelons que nous sommes dans le cas où  $h_1 < h_2 = h^p(\beta, L)$ , et que  $h_1 < h_2 \Rightarrow \sigma_p(h_1) > \sigma_p(h_2)$  si  $\text{pen}_h$  est décroissante en  $h$ . On peut écrire

$$\begin{aligned} \mathbb{P}_f \left\{ |\mathring{f}^{(1)}(y) - \mathring{f}^{(2)}(y)| > C \sigma_p(h_1) \right\} &\leq \mathbb{P}_f \left\{ |\mathring{f}^{(1)}(y) - f(y)| > \frac{C}{2} \sigma_p(h_1) \right\} \\ &\quad + \mathbb{P}_f \left\{ |f(y) - \mathring{f}^{(2)}(y)| > \frac{C}{2} \sigma_p(h_1) \right\} \\ &= \sum_{j=1,2} \mathbb{P}_f \left\{ |\mathring{f}^{(j)}(y) - f(y)| > \frac{C}{2} \sigma_p(h_j) \right\}. \end{aligned}$$

Contrôler la probabilité (1.5.4) revient à contrôler les grandes déviations de l'estimateur  $\mathring{f}^{(j)}$ . Une grande partie de cette thèse est consacrée à la recherche de ce contrôle pour les estimateurs non-linéaires (voir Sections 2.1.3 et 2.2.2). Dans la suite, un exemple de contrôle des grandes déviations est donné pour l'estimateur de la moyenne locale dans le modèle de la régression gaussienne. La décomposition biais-variance joue un rôle particulièrement important, notamment sur le fait que le biais soit croissant et que la variance soit décroissante en  $h$ .

**Grandes déviations : Moyenne locale.** Considérons les observations  $\mathcal{Z}_n$  dans le cadre de la régression gaussienne avec design fixe uniforme comme définie dans la section 1.2.4. Supposons que  $nh^d \in \mathbb{N}^*$  pour simplifier les écritures et définissons l'estimateur de la moyenne locale :

$$\hat{f}^h = \frac{1}{nh^d} \sum_{X_i \in V_h(y)} Y_i, \quad \forall h \in [0, 1].$$

Dans ce cas, l'estimation de la fonction  $f$  revient à approximer celle-ci par une constante sur un voisinage. Ici, nous prendrons  $b = 1$ , si  $\beta \in ]0, 1]$ , les fonctions  $f \in \mathbb{H}(\beta, L)$  sont approchées par des polynômes de degré nul.

On est en présence d'un estimateur *linéaire*. Dans cette exemple la décomposition biais-variance est immédiate.

$$|\hat{f}^h(y) - f(y)| \leq Ldh^\beta + \left| \frac{1}{nh^d} \sum_{X_i \in V_h(y)} \xi_i \right|, \quad \xi_i \sim \mathcal{N}(0, \sigma^2).$$

Supposons ici que la fenêtre  $h \leq h^p(\beta, L)$ . Ainsi, par croissance du biais, décroissance de la variance et compromis biais-variance, on a

$$(1.5.5) \quad Ldh^\beta \leq Ld(h^p(\beta, L))^\beta = \sigma_p(h^p(\beta, L)) \leq \sigma_p(h).$$

On obtient

$$\mathbb{P}_f \left\{ |\hat{f}^h(y) - f(y)| > \frac{C}{2} \sigma_p(h) \right\} \leq \mathbb{P}_f \left\{ \left| \frac{1}{nh^d} \sum_{X_i \in V_h(y)} \xi_i \right| > \frac{C-2}{2} \sigma_p(h) \right\}.$$

Dans le cas gaussien  $\gamma = 2$  ainsi il existe une constante  $c_5 = c_5(\sigma^2) > 0$  telle que

$$\mathbb{P}_f \left\{ \left| \frac{1}{nh^d} \sum_{X_i \in V_h(y)} \xi_i \right| > \frac{C-2}{2} \sigma_p(h) \right\} \leq c_5 \exp \left\{ - \left( \frac{C-2}{2} \text{pen}_h \right)^2 \right\}.$$

Cette inégalité s'obtient en intégrant la probabilité. Si nous choisissons la pénalité  $\text{pen}_h = \sqrt{1 + \ln \frac{h_{max}}{h}}$ , on a

$$\mathbb{P}_f \left\{ |\hat{f}^h(y) - f(y)| > \frac{C}{2} \sigma_p(h) \right\} \leq c_5 \left( \frac{h_{max}}{h} \right)^{-\left(\frac{C-2}{2}\right)^2}.$$

Si l'on choisit  $h, h_{max}$  telles que  $h_{max}h^{-1} \geq n^{c_6}$ , on obtient le résultat avec  $C = 2\sqrt{\frac{2}{c_6}} + 2$ . ■

Nous venons de voir que la méthode de Lepski repose sur des inégalités exponentielles de concentration, qui permettent de contrôler les grandes déviations. Ceci a été une difficulté importante qu'il a fallu surmonter dans cette thèse, notamment pour les estimateurs non-linéaires. Regardons maintenant la règle de sélection générale pour une grille finie d'estimateurs  $\mathring{f}^{(k)}$ .

### Règle de Sélection Générale

Reconsidérons la famille des estimateurs  $\mathcal{F}$ . Nous pouvons maintenant donner la procédure de Lepski établie à partir de la règle de sélection à deux fenêtres. On prend  $f^*(y) = \mathring{f}^{(\hat{k})}(y)$ , où  $\mathring{f}^{(\hat{k})}(y)$  est sélectionné dans  $\mathcal{F}$  suivant la règle suivante :

$$(1.5.6) \quad \hat{k} = \inf \left\{ k = \overline{0, \mathbf{k}_n} : |\mathring{f}^{(k)}(y) - \mathring{f}^{(l)}(y)| \leq C\sigma_p(h_l), \quad l = \overline{k+1, \mathbf{k}_n} \right\},$$

où  $\mathbf{k}_n$  est définie dans (1.5.2),  $\overline{0, \mathbf{k}_n} = 0, \dots, \mathbf{k}_n$  et la constante  $C$  est à choisir comme dans le cas précédent. Pour des exemples, on peut regarder les sections 2.1 et 2.2.  $\mathbf{k}_n$  et  $\mathcal{F}$  sont définis respectivement par (1.5.2) et (1.5.3). Notons que la fenêtre  $h_l$  est une fonction décroissante en  $l$  et la variance pénalisée  $\sigma_p(h_l)$  est croissante en  $l$ .

### Interprétation de la règle de Lepski

La règle de Lepski peut être expliquée par la description suivante. Considérons les fenêtres  $0 < h_{\min} \leq h_{\mathbf{k}_n} < \dots < h_0 = h_{\max}$ , on en déduit l'inclusion  $\emptyset \neq V_{h_{\mathbf{k}_n}}(y) \subset \dots \subset V_{h_1}(y)$ . Dans ce cas, la variance de notre estimateur  $\mathring{f}^{(k)}$  est croissante  $\text{Var}(\mathring{f}^{(k)}) \leq \text{Var}(\mathring{f}^{(k+1)})$ ,  $k = 1, \dots, \mathbf{k}_n - 1$ . En d'autres termes, le biais est en général décroissant quand le voisinage diminue  $V_{h_k}(y)$ . La méthode de Lepski peut être interprétée comme une procédure de tests multiples où l'hypothèse, que  $f$  est un polynôme sur  $V_{h_k}(y)$ , qui est testé contre le fait que l'écart entre les estimateurs est important. Notons cette hypothèse  $H_0(k) : f_\theta|_{V_{h_k}(y)} = f|_{V_{h_k}(y)}$ , où  $\theta$  est la suite des coefficients des polynômes de Taylor.

Supposons  $H_0(\mathbf{k}_n)$  est vraie, nous testons successivement si  $H_0(k-1)$  est acceptable à condition que l'hypothèse  $H_0(l)$  soit vérifiée pour tout  $l \geq k$ . Une fois le test de  $H_0(k-1)$  rejeté, nous choisissons l'estimateur  $\mathring{f}^{(k)}$  correspondant à la dernière hypothèse. Le point principal est donc de bien définir les tests d'acceptation pour  $H_0(k-1)$  : la méthode de Lepski accepte  $H_0(k-1)$  si l'événement  $|\mathring{f}^{(k-1)}(y) - \mathring{f}^{(l)}(y)| \leq CS_n(l)$  est vrai pour tout  $l \geq k-1$  avec des valeurs de seuillage  $CS_n(l)$  appropriées. Cette quantité  $CS_n(l)$  peut être vu comme l'écart-type pénalisé de  $\mathring{f}^{(l)}$  à une constante  $C$ -près.

### Implémentation de la règle de Lepski

Soit  $A_{j,l} = \left\{ |\mathring{f}^{(j)}(y) - \mathring{f}^{(l)}(y)| \leq C\sigma_p(h_l) \right\}$ . Pour l'implémentation, on constate qu'il est nécessaire de calculer tous les estimateurs  $\mathring{f}^{(l)}$  de la grille  $\mathcal{F}$ . Ceci donne la règle suivante :

on dit qu'un estimateur  $\hat{f}^{(j)}$  est *admissible* si l'événement  $\bigcap_{j < l \leq \mathbf{k}_n} A_{j,l}$  est réalisé.

1. On commence par vérifier si l'estimateur  $\hat{f}^{(0)}$  est admissible, ,
2. si il est admissible, on choisit  $f^* = \hat{f}^{(0)}$ , sinon on passe au suivant  $\hat{f}^{(1)}$ .
3. on s'arrête au premier instant où l'estimateur  $\hat{f}^{(j)}$  est admissible et tous les précédents  $\hat{f}^{(0)}, \dots, \hat{f}^{(j-1)}$  ne le sont pas. Ainsi, on prend  $f^* = \hat{f}^{(j)}$  où  $j = \hat{k}$  est la plus petite des valeurs admissibles.

Donc, dès le début de l'algorithme, pour vérifier que  $\hat{f}^{(0)}$  est admissible ou non, on construit tous les estimateurs de la grille  $\mathcal{F}$ .

Pour contourner ce problème, on peut utiliser l'implémentation de [Katkovnik \[1999\]](#) avec les ICI (en anglais : *Intersection of Confidence Intervals*), qui ne nécessite pas de calculer tous les estimateurs  $\hat{f}^{(l)}$  de la grille  $\mathcal{F}$ . La règle ICI est la suivante :

### Implémentation des ICI

1. Classons par ordre croissant les fenêtres  $h_{\mathbf{k}_n} < h_{\mathbf{k}_n-1} < \dots < h_1 < h_0$ .
2. Considérons les intersections des intervalles de confiance  $\mathcal{I}_j = \bigcap_{i=\mathbf{k}_n}^j \mathcal{D}_i$  pour  $j \in \{\mathbf{k}_n, \mathbf{k}_n - 1, \dots, 1, 0\}$ , avec

$$\mathcal{D}_i = \left[ \hat{f}^{(i)}(y) - C_{\varsigma_j}, \hat{f}^{(i)}(y) + C_{\varsigma_j} \right],$$

$\varsigma_j = \text{Std}(\hat{f}^{(j)}(y))$  est l'écart-type de  $\hat{f}^{(j)}(y)$ .

3. Soit  $j^+$  l'entier tel que pour tout  $j > j^+$  on a  $\mathcal{I}_j \neq \emptyset$ ,  $\mathcal{I}_{j^+} \neq \emptyset$  et  $\mathcal{I}_{j^+-1} = \emptyset$ . Ainsi, l'estimateur adaptatif est  $\hat{f}^{(j^+)}(y)$ .

Une illustration des ICI est donnée par la figure 1.10. On constate que si  $j$  diminue, alors l'écart-type associé  $\varsigma_j$  diminue. Ainsi, les intervalles de confiance diminuent quand  $j$  diminue. L'idée générale, cachée derrière les ICI, est que l'estimation du biais est plus petite, à constante près, que la variance tant que les intersections sont non-vides. La règle s'arrête quand l'estimation du biais et la variance sont proches.

La figure 1.11 permet de constater quelle fenêtre la méthode a choisie. La fenêtre  $h_{\mathbf{k}-1}$  est trop petit, elle n'est pas optimale, nous pouvons encore l'agrandir. On passe à la fenêtre supérieure, on voit que celle-ci convient. La fenêtre  $h_{\mathbf{k}-3}$  est trop grande et implique que  $\mathcal{I}_{\mathbf{k}-3}$  est vide.

Si  $\hat{f}^{(j)}$  est un estimateur à noyau avec un noyau  $K_h$ , l'écart-type  $\varsigma_j$  peut facilement être calculé par  $\varsigma \|K_{h_j}\|_2$  (norme  $\ell_2$  du noyau,  $\varsigma$  est l'écart-type du bruit). Dans les autres cas, on peut remplacer  $\varsigma_j$  par  $\sigma_p(h_j)$ .

Asymptotiquement, on peut démontrer que l'estimateur ICI-adaptatif est proche de  $f^*$ ,  $\hat{f}^{(j^+)}(y) \approx f^*(y)$  (voir [Goldenshluger et Nemirovski \[1997\]](#)). Pour les calculs théoriques, nous utiliserons l'estimateur  $f^*$  choisi par la règle de Lepski. En pratique, nous utiliserons  $\hat{f}^{(j^+)}$  l'estimateur choisi par la règle des ICI.



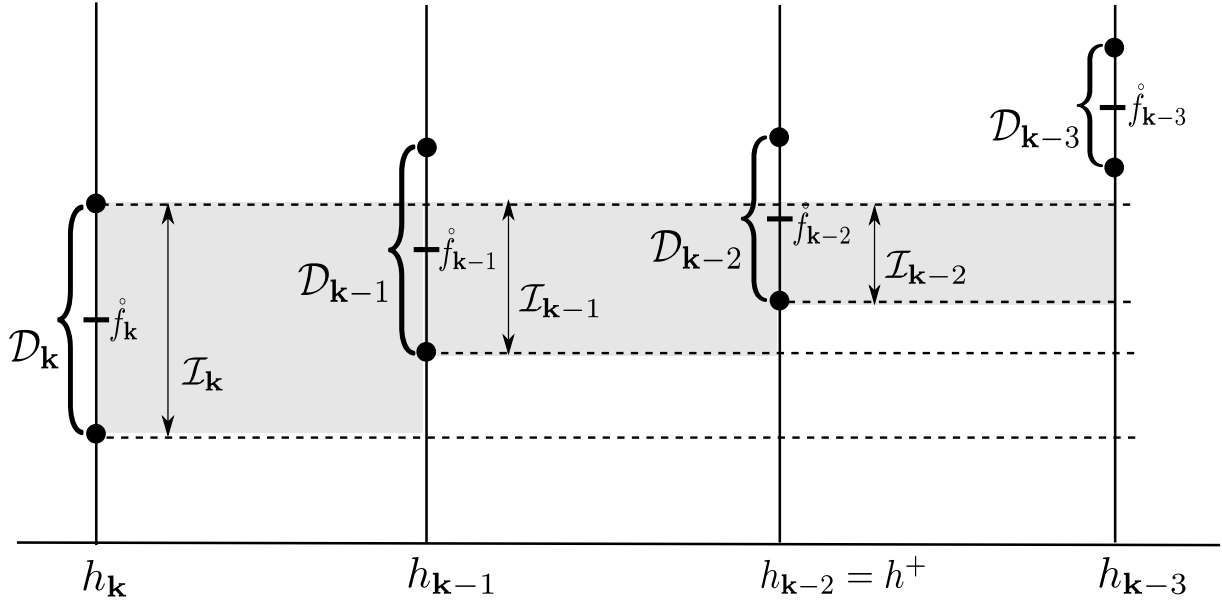


FIGURE 1.10 – Illustration de la règle ICI développée par [Katkovnik \[1999\]](#) ( $\mathring{f}_j = \mathring{f}^{(j)}(y)$  et  $h^+ = h_{j^+}$ ).

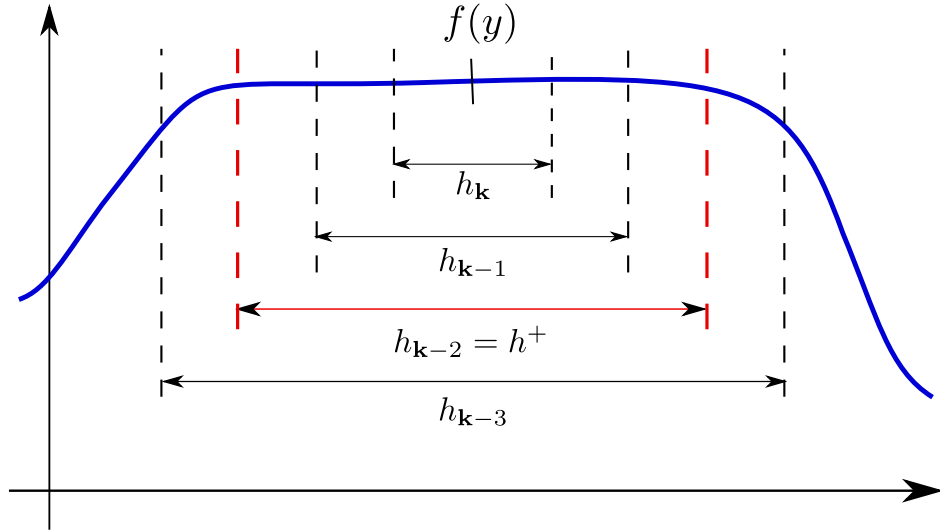


FIGURE 1.11 – Choix de la fenêtre  $h^+ = h_{j^+}$  en traçant la fonction cible (exemple où la fonction est localement constante).

Comme nous venons de le voir, l'interprétation de la méthode de Lepski par les tests statistiques et son implémentation très simple par la méthode ICI ont fait le succès de cette méthode.

Nous pouvons démontrer que l'estimateur  $f^*$  est adaptatif.

**Proposition 2.** *Pour tout  $\beta, L \in [0, b] \times \mathbb{R}_+$  on a*

$$\sup_{f \in \mathbb{H}_d(\beta, L)} \mathbb{E}_f |f^*(y) - f(y)| \leq (Cc_1^* + c_2^*)\psi_n^p(\beta) + T_n,$$

où  $c_1^*, c_2^*$  sont des constantes positives et  $T_n$  est un terme de reste que l'on définira dans la preuve.

Dans la procédure (1.5.6), on peut remplacer  $C$  par une constante  $C'$  telle que  $C' > C$ . Le résultat du Proposition 2 reste vrai en remplaçant  $C$  par  $C'$  dans la borne supérieure (voir Preuve).

**Idée de la preuve.** Soit l'entier  $\kappa$  défini de la façon suivante.

$$2^{-\kappa} h_{\max} \leq h^p(\beta, L) < 2^{-\kappa+1} h_{\max}.$$

On a aussi

$$\begin{aligned} \mathbb{E}_f |f^{\circ(\hat{k})}(y) - f(y)| &\leq \mathbb{E}_f |f^{\circ(\hat{k})}(y) - f(y)| \mathbb{I}_{\hat{k} \leq \kappa} + \mathbb{E}_f |f^{\circ(\hat{k})}(y) - f(y)| \mathbb{I}_{\hat{k} > \kappa} \\ (1.5.7) \quad &:= R_1(f) + R_2(f). \end{aligned}$$

Premièrement, nous contrôlons  $R_1$ . Ainsi

$$|f^{\circ(\hat{k})}(y) - f(y)| \leq |f^{\circ(\hat{k})}(y) - f^{\circ(\kappa)}(y)| + |f^{\circ(\kappa)}(y) - f(y)|.$$

La définition de  $\hat{k}$  entraîne que

$$|f^{\circ(\hat{k})}(y) - f^{\circ(\kappa)}(y)| \mathbb{I}_{\hat{k} \leq \kappa} \leq C\sigma_p(h_\kappa).$$

Par définition de  $\kappa$ , l'estimateur  $f^{\circ(\kappa)}$  est le plus proche de l'estimateur minimax pénalisé dans la grille  $\mathcal{F}$ . Il est facile de démontrer qu'il existe une constante  $c_1^*$  telle que  $\sigma_p(h_\kappa) \leq c_1^* \psi^p(\beta)$  et une constante  $c_2^*$  telle que

$$\mathbb{E}_f |f^{\circ(\kappa)}(y) - f(y)| \leq c_2^* \psi^p(\beta).$$

D'où  $R_1(f) \leq (Cc_1^* + c_2^*)\psi^p(\beta)$ .

Maintenant, il nous faut borner  $R_2$ . Appliquons l'inégalité de Cauchy-Schwarz

$$\begin{aligned} R_2(f) &= \sum_{k > \kappa}^{\mathbf{k}_n} \mathbb{E}_f |f^{\circ(k)}(y) - f(y)| I_{[\hat{k}=k, G]} \\ &\leq \sum_{k > \kappa} (\mathbb{E}_f |f^{\circ(k)}(y) - f(y)|^2)^{1/2} \sqrt{\mathbb{P}_f \{\hat{k} = k\}}. \end{aligned}$$

Remarquons que l'indice  $\hat{k}$  est aléatoire, nous le contrôlons à l'aide de la somme qui rend l'indexation déterministe. Cela va nous permettre d'utiliser un contrôle des grandes déviations pour l'estimateur en question. Car si le paramètre  $l$  de l'estimateur localement polynômial  $\hat{f}^{(k)}$  est aléatoire (choisi à partir des données), alors cet estimateur n'est plus localement paramétrique. Ainsi le contrôle des grandes déviations est plus difficile à obtenir. Par exemple, les inégalités de concentration dit de *Bernstein* ou de *Berstein* (que l'on peut trouver dans [Boucheron, Bousquet, et Lugosi \[2004\]](#)) peuvent être utilisées pour les estimateurs linéaires seulement si la fenêtre est déterministe. Notons que pour tout  $k \geq \kappa + 1$  et par définition de  $\hat{k}$  (1.5.6)

$$\{\hat{k} = k\} = \cup_{l \geq k} \left\{ |\hat{f}^{(k-1)}(y) - \hat{f}^{(l)}(y)| > C\sigma_p(h_l) \right\}.$$

Notons que  $\sigma_p(h_l)$  est strictement croissant en  $l$ , donc,

$$\begin{aligned} \{\hat{k} = k\} &\subseteq \left\{ |\hat{f}^{(k-1)}(y) - f(y)| > 2^{-1}C\sigma_p(h_{k-1}) \right\} \\ &\cup \left[ \cup_{l \geq k} \left\{ |\hat{f}^{(l)}(y) - f(y)| > 2^{-1}C\sigma_p(h_l) \right\} \right]. \end{aligned}$$

On arrive à l'inégalité suivante : pour tout  $k \geq \kappa + 1$

$$\begin{aligned} \mathbb{P}_f(\hat{k} = k) &\leq \mathbb{P}_f \left\{ |\hat{f}^{(k-1)}(y) - f(y)| > 2^{-1}C\sigma_p(h_{k-1}) \right\} \\ &+ \sum_{l \geq k} \mathbb{P}_f \left\{ |\hat{f}^{(l)}(y) - f(y)| > 2^{-1}C\sigma_p(h_l) \right\}. \end{aligned}$$

A ce moment de la preuve, nous pouvons donner une hypothèse pour finir la démonstration.

**Hypothèses 2.** Pour tout  $k \geq \kappa$ , il existe une constante  $c_3^*$  telle que

$$\mathbb{P}_f \left\{ |\hat{f}^{(k)}(y) - f(y)| > 2^{-1}C\sigma_p(h_k) \right\} \leq c_3^* \exp \left\{ - \left( \frac{C - c_4^*}{2} \text{pen}_{h_k} \right)^\gamma \right\}.$$

Dans cette thèse, nous prendrons la pénalité  $\text{pen}_{h_k} = (1 + \ln [h_{\max} h_k^{-1}])^{1/\gamma} = (1 + k \ln 2)^{1/\gamma}$ . Avec l'hypothèse précédente, nous contrôlons

$$\mathbb{P}_f(\hat{k} = k) \leq c_5^* \exp \left\{ - \left( \frac{C - c_4^*}{2} \right)^\gamma (k - 1) \ln 2 \right\},$$

où  $c_5^* = c_3^* \sum_{l=0}^{+\infty} \exp \left\{ - \left( \frac{C - c_4^*}{2} \right)^\gamma l \ln 2 \right\}$ . Ainsi, il est facile de contrôler le risque suivant. Il existe une constante  $c_6^*$  telle que

$$(\mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^2)^{1/2} \leq c_6^* \sigma_p(h_k).$$

On obtient ce résultat en intégrant la probabilité de l'hypothèse 2. D'où

$$R_2(f) \leq \sum_{k > \kappa} c_6^* \sigma_p(h_k) \sqrt{c_5^*} \exp \left\{ -\frac{1}{2} \left( \frac{C - c_4^*}{2} \right)^\gamma (k - 1) \ln 2 \right\}.$$

Avec le choix de  $C = 2(2d)^{1/\gamma} + c_4^*$ , on obtient

$$R_2(f) \leq T_n = \frac{1}{(nh_{\max}^d)^{1/\gamma}} \sum_{k > 0} c_6^* \sqrt{c_5^*} (1 + k \ln 2)^{1/\gamma} \exp \left\{ -\frac{1}{2} \left( \frac{C - c_4^*}{2} \right)^\gamma (k - 1) \ln 2 \right\}.$$

Avec cette inégalité, l'assertion du proposition est démontrée.

La recherche, du contrôle des grandes déviations, est une partie technique mais indispensable pour démontrer les résultats d'adaptation. Cette thèse consacre une partie importante à la recherche de ce type d'inégalités notamment pour les estimateurs non-linéaires. De nombreux papiers sont consacrés à ces résultats dans le but de l'adaptation (voir notamment Talagrand [1994], Ledoux [1997], Massart [2000, 2007], Boucheron, Bousquet, et Lugosi [2004] et Goldenshluger et Lepski [2009b]).

Nous insistons sur le fait que le biais et la variance de notre estimateur doivent être respectivement croissant et décroissant (voir l'équation (1.5.5)). Pour les fonctions anisotropes, le biais n'est pas monotone. Dans ce sens, les articles de Kerkycharian, Lepski, et Picard [2001], Klutchnikoff [2005], Goldenshluger et Lepski [2008] et Goldenshluger et Lepski [2009a] utilisent une méthode s'appuyant sur une comparaison des biais. La notion d'emboîtement est inévitable. En effet il est nécessaire de trouver un ordre sur l'ensemble des fenêtres pour que la variance soit monotone. En général, pour l'approche localement paramétrique, le nombre d'observations appartenant à  $V_h(y)$  augmente quand la fenêtre augmente, d'où la diminution de la variance.

Dans cette procédure, nous avons un "point d'appui",  $\hat{f}^{(\kappa)}$  l'estimateur minimax pénalisé, qui est notre *oracle*. Ainsi, la méthode de Lepski peut donner des inégalités oracle (voir Goldenshluger et Lepski [2009a]), mais cela demande des résultats un peu plus techniques (bornes uniformes de fonctions aléatoires, voir Goldenshluger et Lepski [2009b]).

La méthode de Lepski (Lepski, Mammen, et Spokoiny [1997]) nous a permis de construire des estimateurs *bayésien et de Huber adaptatifs optimaux*. Nous présentons dans le chapitre suivant les résultats obtenus dans cette thèse notamment grâce à cette méthode. Dans la section 2.3, la performance de cette méthode est évaluée (par rapport à l'oracle). Une illustration de la calibration de la constante  $C$  est aussi présentée.



# Chapitre 2

## Résultats et Perspectives

Après avoir défini les modèles, les nouvelles méthodes d'estimation et la procédure adaptative, nous présentons les différents théorèmes qui mesurent les performances des estimateurs proposés. D'un point de vue minimax dans un premier temps, nous développons ensuite des procédures adaptatives basées sur l'idée de Lepski. Comme nous l'avons déjà souligné ci-dessus, nos résultats sont déduits du contrôle des grandes déviations de nos différents estimateurs. Ces inégalités font l'objet de sous-sections pour chaque méthode d'estimation, et peuvent être étudiées de façon indépendante des autres résultats. En fin de paragraphe, nous parlerons des suites à donner à ces travaux, ainsi que des problèmes ouverts.

### 2.1 Approche Bayésienne

Nous donnons les résultats obtenus pour l'estimateur bayésien dans le modèle de la régression générale. Tout d'abord nous précisons quelques hypothèses : la vitesse de convergence minimax est supposée connue  $\varphi_{n,\gamma}(\beta) = n^{-\frac{\beta}{\gamma\beta+d}}$  sur les espaces de Hölder  $\mathbb{H}_d(\beta, L)$ , où  $1 \leq \gamma \leq 2$  est une constante qui dépend de la densité  $g(.,.)$  des observations. Une propriété remarquable de l'estimateur bayésien est qu'il s'adapte à la densité du modèle pour être optimal au sens des vitesses de convergence minimax. Comme tout estimateur localement paramétrique, le choix de la fenêtre est primordial pour construire un estimateur qui ne dépend pas de la fonction à estimer. Nous proposons une procédure adaptative pour le choix de la fenêtre, basée sur l'idée de Lepski, qui construit un estimateur adaptatif optimal. Nous verrons quelques exemples d'optimalité pour certains modèles (section 4.3). Comme nous l'avons vu dans la section 1.5.2, le contrôle des grandes déviations pour l'estimateur bayésien indexé par  $h$  doit être obtenu pour tout  $h$ . Ici, l'estimateur bayésien n'est pas linéaire et on ne connaît pas sa forme explicite. Nous utiliserons le fait qu'il minimise un critère particulier pour contrôler les grandes déviations. Une fois les résultats énoncés pour le modèle de régression générale, nous étudions quelques exemples comme la régression gaussienne et la régression inhomogène de Poisson avec la vitesse  $\varphi_{n,2}(\beta)$ , la régression  $\alpha$  avec  $\varphi_{n,1+2\alpha}(\beta)$  et

la régression multiplicative uniforme avec  $\varphi_{n,1}(\beta)$ .

Nous proposons une procédure adaptative pour l'estimateur bayésien. Cette procédure permet un choix à partir des observations de la fenêtre d'estimation (en anglais : bandwidth). Les résultats obtenus sont de type minimax adaptatifs pour le risque ponctuel. En particulier notre estimateur adaptatif atteint la vitesse  $\varphi_{n,\gamma}(\beta)$ ,  $\forall \beta \in ]0, b]$ , nous verrons plusieurs modèles dans lesquels cette vitesse est optimale. Notons toutefois que l'adaptation pour l'approche bayésienne est déjà apparue, par exemple dans les travaux de Ghosal, Lember, et Van der Vaart [2008] et Van der Vaart et Van Zanten [2009].

### 2.1.1 Recherche de la Vitesse Minimax

Rappelons le cadre dans lequel on travaille.

1. Le modèle est celui de la régression générale défini dans (1.2.2) avec une densité  $g(\cdot, \cdot)$ ,
2.  $f \in \mathbb{H}_d(\beta, L, M)$  (Définition 1),
3. on utilise l'estimateur bayésien  $\hat{f}$  défini dans (1.3.7), avec la constante  $m$  à choisir qui correspond à l'inverse de la puissance du rapport de vraisemblance,
4. on suppose que la vitesse minimax est connue et égale à  $\varphi_{n,\gamma}(\beta) = n^{-\frac{\beta}{\gamma\beta+d}}$  avec  $\gamma$  connue.

Pour obtenir la borne supérieure du risque maximal (Définition 6), nous avons besoin des hypothèses suivantes. Pour cela, on introduit les notations suivantes. Pour tout  $\theta \in \Theta(M)$ , soit  $U_n = N_h [\Theta(M) - \theta]$ . Donnons aussi  $N_h = (nh^d)^{1/\gamma}$  où  $\gamma \geq 1$  et  $\forall u \in U_n$  le pseudo rapport de vraisemblance

$$(2.1.1) \quad Z_{n,\theta}(u) = \frac{L_h(\theta + u N_h^{-1}, \mathcal{Z}_n)}{L_h(\theta, \mathcal{Z}_n)}.$$

Soit  $\mathcal{H}_n, n > 1$  le sous-intervalle suivant de  $(0, 1)$ .

$$(2.1.2) \quad \mathcal{H}_n = [h_{\min}, h_{\max}], \quad h_{\min} = (\ln n)^{\frac{1}{\gamma d + d^2}} n^{-1/d}, \quad h_{\max} = (\ln n)^{-\frac{1}{\gamma b + d}}.$$

Dans la suite, nous considérons seulement les valeurs de  $h$  appartenant à  $\mathcal{H}_n$ . Soit  $f_\theta(x)$ , avec  $\theta \in \Theta(M)$ , le polynôme local d'approximation de  $f$  sur  $V_h(y)$  et soit  $b_h$  l'erreur d'approximation correspondante, i.e.

$$(2.1.3) \quad b_h = \sup_{x \in V_h(y)} |f_\theta(x) - f(x)|.$$

Finalement, définissons  $\mathcal{N}(h) = (b_h \times N_h)^\gamma$  et

$$(2.1.4) \quad \mathcal{E}_h = \exp \{ \mathcal{N}(h) \}.$$

Soit le sous-intervalle  $\Gamma_\delta \subseteq \Theta(M)$  tel que  $\int_{\Gamma_\delta} du = \delta^{D_b}$ ,  $\delta > 0$ .

**Hypothèses 3.** On suppose qu'il existe plusieurs constantes  $\tau, m > 0$  et  $c_1, c_2, c_3, C_1, C_2, s_1, s_2 > 0$  telles que pour tout  $n > 1$ ,  $f \in \mathbb{H}_d(\beta, L, M)$ ,  $\theta \in \Theta(M)$  et  $h \in \mathcal{H}_n$ , on a

1.  $\int_{\Gamma_\delta} \mathbb{E}_f \left[ 1 - Z_{n,\theta}^{1/m}(u) \right]_+ du \leq C_1 \mathcal{E}_h^{c_1} \delta^{D_b + \tau}, \quad \forall \delta < s_1,$
2.  $\mathbb{E}_f Z_{n,\theta}^{1/m}(u) \leq C_2 \mathcal{E}_h^{c_2} \exp \{ -c_3 \|u\|_1^\gamma \}, \quad \forall u \in U_n : \|u\|_1 \geq s_2 \mathcal{N}(h).$

où  $a_+ = \max(a, 0)$ ,  $a \in \mathbb{R}$ .

**Remarque 3.** On suppose aussi connues les constantes  $c_3, s_2, \tau, m$  and  $\gamma$ , l'existence des autres constantes est suffisante. Notons que le processus  $Z_{n,\theta}(\cdot)$  n'est pas toujours défini, par exemple si la mesure de probabilité  $\mathbb{P}_{f_{\theta+uN_h^{-1}}}$  n'est pas absolument continue par rapport à  $\mathbb{P}_{f_\theta}$ . Dans ce cas, le problème est ouvert.

Par le passé, [Has'minskii et Ibragimov \[1981\]](#) ont énoncé des hypothèses similaires pour les modèles paramétriques. Ils utilisent l'estimateur bayésien  $\int \|t - u\| L(u) du$  mais sans la puissance  $1/m$ . Ainsi, les hypothèses ci-dessus sont plus faibles que celles de [Has'minskii et Ibragimov \[1981\]](#) et [Chichignoud \[2010a\]](#) (voir Chapitre 4). Ces hypothèses sont données par [Chichignoud \[2010b\]](#) (Chapitre 3). La continuité du processus  $Z_{n,\theta}(\cdot)$ , et donc de la vraisemblance, n'est pas nécessaire pour l'approche bayésienne. En revanche,  $Z_{n,\theta}$  est défini comme le rapport de deux mesures de probabilités, il est donc nécessaire que  $\mathbb{P}_{f_{\theta+uN_h^{-1}}}$  soit absolument continue par rapport à  $\mathbb{P}_{f_\theta}$ . Le prochain théorème montre qu'avec un choix judicieux de la fenêtre  $h$ , l'estimateur bayésien local est optimal d'un point de vue minimax sur les espaces de Hölder. Posons  $\bar{h} = (L^\gamma n)^{-\frac{1}{\gamma\beta+d}}$  et soit  $\bar{f}^h(y) = \hat{\theta}_{0,\dots,0}(\bar{h})$  donné par (1.3.4), (1.3.6) et (1.3.7) avec  $h = \bar{h}$ .

**Théorème 1.** Soient  $\beta > 0$ ,  $L > 0$  et  $M > 0$  fixées. Supposons que les hypothèses 3 sont vérifiées. Alors, il existe une constante  $C_*$  telle que pour tout  $n \in \mathbb{N}^*$  satisfaisant  $N_{\bar{h}} \geq 1$ ,

$$\varphi_{n,\gamma}^{-q}(\beta) R_{n,q} \left[ \bar{f}^h(y), \mathbb{H}_d(\beta, L, M) \right] \leq C_*, \quad \forall q \geq 1.$$

**Remarque 4.** Sous les hypothèses 3, nous savons construire un estimateur minimax. Nous verrons dans les sections 2.1.4 et 3.4 dans quels modèles ces hypothèses sont vérifiées. Les estimateurs de type bayésien ou de Pitman ont été introduits dans le livre de [Has'minskii et Ibragimov \[1981\]](#) pour l'estimation paramétrique, dans le but d'atteindre la vitesse optimale pour différents modèles où la vitesse varie. Dans la preuve (section 3.5), la constante  $C_*$  est exprimée de façon explicite. Notons que la borne inférieure n'est pas calculée car nous avons supposé la vitesse minimax connue. Dans les exemples (section 2.1.4), les bornes inférieures du risque minimax sont calculées.



### 2.1.2 Procédure Adaptative

Cette section est vouée à l'estimation adaptative sur la collection de classes  $\left\{ \mathbb{H}_d(\beta, L, M) \right\}_{\beta, L, M}$ . Nous n'allons pas imposer de restriction sur les valeurs possibles de  $L, M$ . En effet,  $L, M$  joue un rôle secondaire en intervenant dans les constantes, à l'inverse de  $\beta$  qui a une influence sur la vitesse. Nous supposons que  $\beta \in (0, b]$ , où  $b$ , est un entier choisi arbitrairement *a priori*. Notons aussi que la constante  $M$  peut être estimée par un estimateur non nécessairement optimal (voir Chichignoud [2010a] ou Chapitre 4, Section 4.3). En revanche, il n'existe aucune méthode capable d'estimer  $\beta$  et  $L$ .

Soit  $\Phi$  la famille de normalisations suivante :

$$\phi_{n,\gamma}(\beta) = \left( \frac{\rho_{n,\gamma}(\beta)}{n} \right)^{\frac{\beta}{\gamma\beta+d}}, \quad \rho_{n,\gamma}(\beta) = \left[ 1 + \frac{\gamma(b-\beta)}{(\gamma b+d)(\gamma\beta+d)} \ln n \right]^{\frac{1}{\gamma}}, \quad \beta \in (0, b].$$

On remarque que  $\phi_{n,\gamma}(b) = \varphi_{n,\gamma}(b)$  et  $\rho_{n,\gamma}(\beta) \sim (\ln n)^{\frac{1}{\gamma}}$  pour tout  $\beta \neq b$ . On constate une perte dans la vitesse adaptative. Ceci est en général toujours vrai pour l'estimation ponctuelle (voir notamment Lepski [1990], Tsybakov [1998], Klutchnikoff [2005] et Chichignoud [2010a]). Nous verrons que cette perte est nécessaire pour contrôler les grandes déviations pour l'estimateur bayésien local.

**Construction de l'estimateur  $\Phi$ -adaptatif.** Comme mentionné dans la section 1.5.2, la construction de notre procédure d'estimation comprend deux étapes. Premièrement, nous définissons la famille des estimateurs bayésiens locaux. Ensuite, s'appuyant sur la méthode de Lepski, nous proposons une sélection adaptative pour cette famille.

1ère étape : Collection d'estimateurs bayésiens locaux. Posons

$$h_k = 2^{-k} h_{\max}, \quad k = 0, \dots, \mathbf{k}_n,$$

où  $\mathbf{k}_n$  est le plus grand entier tel que  $h_{\mathbf{k}_n} \in \mathcal{H}_n$  défini dans (2.1.2). Soit

$$\hat{\mathcal{F}} = \left\{ \hat{f}^{(k)}(y) = \hat{\theta}_{0,\dots,0}(h_k), \quad k = 0, \dots, \mathbf{k}_n \right\},$$

où  $\hat{\theta}_{0,\dots,0}(h_k)$  est donné par (1.3.4), (1.3.6) et (1.3.7) avec  $h = h_k$ .

2ème étape : Sélection parmi la collection  $\hat{\mathcal{F}}$ . On pose  $\hat{f}^*(y) = \hat{f}^{(\hat{k})}(y)$ , où  $\hat{f}^{(\hat{k})}(y)$  est sélectionné dans  $\hat{\mathcal{F}}$  avec la règle de décision suivante :

$$(2.1.5) \quad \hat{k} = \inf \left\{ k = \overline{0, \mathbf{k}_n} : |\hat{f}^{(k)}(y) - \hat{f}^{(l)}(y)| \leq C S_n(l), \quad l = \overline{k+1, \mathbf{k}_n} \right\}.$$

Ici, on a utilisé les notations suivantes, pour  $q \geq 1$  :

$$(2.1.6) \quad C = \left( s_2^\gamma \vee \frac{2^{2\gamma+2}(\gamma+dq)}{c_3\gamma(1 \wedge \tau D_b^{-1})} \right)^{\frac{1}{\gamma}}, \quad S_n(l) = \left[ \frac{1+l \ln 2}{n(h_l)^d} \right]^{\frac{1}{\gamma}}, \quad l = 0, 1, \dots, \mathbf{k}_n.$$

**Théorème 2.** Soit  $b > 0$  fixé. Supposons que les hypothèses 3 soient vérifiées, alors pour tout  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$  et  $q \geq 1$

$$\limsup_{n \rightarrow \infty} \phi_{n,\gamma}^{-q}(\beta) R_{n,q} \left[ \hat{f}^*(y), \mathbb{H}_d(\beta, L, M) \right] < \infty.$$

**Remarque 5.** Ce théorème correspond au théorème 5 dans le chapitre 3. L’assertion de ce théorème montre que l’estimateur proposé  $\hat{f}^*(y)$  est  $\Phi$ -adaptatif. Cela implique que la famille de normalisations  $\Phi$  est admissible. Par exemple, nous pouvons démontrer l’optimalité de  $\Phi$  dans la régression multiplicative uniforme avec le critère de Klutchnikoff [2005] (voir Section 4.6). Nous avons choisi d’écrire le théorème de façon asymptotique pour simplifier la présentation, au vue de la preuve (section 3.5) nous pourrions écrire la borne supérieure sous la forme  $\phi_{n,\gamma} + R_n$  comme présentée dans la section 1.5.2. La connaissance des constantes  $s_2$ ,  $\gamma$ ,  $c_3$  est nécessaire pour construire la constante  $C$  intervenant dans la procédure. En pratique, on peut calibrer la constante  $C$  sur des jeux de données et ainsi éviter ce problème.

### 2.1.3 Grandes Déviations

Nous avons précisé dans la section 1.5.2 que le contrôle des grandes déviations est indispensable pour obtenir l’adaptation. Ainsi, nous avons pu établir un nouveau type d’inégalités exponentielles. Il existe  $\varepsilon_0$ ,  $\omega_1$ ,  $\omega_2$  et  $\omega_3$  des constantes positives. Pour tout  $n > 1$ ,  $h \in \mathcal{H}_n$  et  $f \in \mathbb{H}_d(\beta, L, M)$ , alors  $\forall \varepsilon \geq \varepsilon_0 \times \mathcal{N}(h)$

$$\mathbb{P}_f \left( N_h | \hat{f}^h(y) - f(y) | \geq \varepsilon \right) \leq \omega_1 \mathcal{E}_h^{\omega_2} \exp \{ -\omega_3 \varepsilon^\gamma \},$$

où l’estimateur bayésien  $\hat{f}^h(y)$  est donné par (1.3.4), (1.3.6) et (1.3.7). Ce résultat est énoncé dans la proposition 3, la preuve est disponible dans la section 3.5.4. La nouveauté ici est la puissance  $\gamma$  dans l’exponentielle. En effet la puissance conditionne la vitesse  $\phi_{n,\gamma}(\beta)$  (voir aussi Has’minskii et Ibragimov [1981]). A noter aussi l’apparition du terme  $\mathcal{E}_h = \exp\{b_h N_h\}$ . Ce terme peut éventuellement tendre vers l’infini si le nombre d’observations augmente. Pour le contrôler, nous utiliserons le compromis biais-variance, i.e.  $\mathcal{N}(h) = b_h N_h \leq L d h^\beta \times (n h^d)^{1/\gamma}$ , nous pouvons voir ce terme comme le rapport biais/variance. Pour l’estimateur minimax, la borne supérieure du risque s’obtient par intégration de cette probabilité.

Dans la suite, nous donnons quelques idées essentielles de la preuve. Notamment comment utiliser les hypothèses 3 et le critère bayésien.

**Idée de la preuve.** La définition de  $\hat{\theta}(h)$  et  $\theta = \theta(f, y, h)$  les coefficients du polynôme de Taylor associé à  $f$ , impliquent que  $\forall \varepsilon > 0$

$$\begin{aligned} \mathbb{P}_f \left( N_h | \hat{f}^h(y) - f(y) | \geq \varepsilon \right) &\leq \mathbb{P}_f \left( N_h | \hat{\theta}_0(h) - \theta_0 | \geq \varepsilon \right) \\ (2.1.7) \qquad \qquad \qquad &\leq \mathbb{P}_f \left( N_h \| \hat{\theta}(h) - \theta \|_1 \geq \varepsilon \right). \end{aligned}$$

Quelques remarques sont à faire. Premièrement, par définition on a  $\theta \in \Theta(M)$ . Rappelons que  $\hat{\theta}(h)$  minimise  $\pi_h$  défini par (1.3.6) et donc, l'inclusion suivante est vraie si  $\hat{\theta}(h) \in \Theta(M)$ .

$$(2.1.8) \quad \left\{ N_h \|\hat{\theta}(h) - \theta\|_1 \geq \varepsilon \right\} \subseteq \left\{ \inf_{N_h \|t - \theta\|_1 \geq \varepsilon} \pi_h(t) \leq \pi_h(\theta) \right\}.$$

En outre,

$$\begin{aligned} \pi_h(t) &= N_h^{-1} \int_{\Theta} \|N_h(t - u)\|_1 [L_h(u, \mathcal{Z}_n)]^{1/m} du \\ &= N_h^{-D_b-1} \int_{U_n} \|N_h(t - \theta) - u\|_1 [L_h(\theta + uN_h^{-1}, \mathcal{Z}_n)]^{1/m} du \\ &= N_h^{-D_b-1} [L_h(\theta, \mathcal{Z}_n)]^{1/m} \int_{U_n} \|N_h(t - \theta) - u\|_1 Z_{n,\theta}^{1/m}(u) du. \end{aligned}$$

D'où,  $\tau_n = N_h(\hat{\theta}(h) - \theta)$  est le minimiseur de

$$\chi_n(s) = \int_{U_n} \|s - u\|_1 \frac{Z_{n,\theta}^{1/m}(u)}{\int_{U_n} Z_{n,\theta}^{1/m}(v) dv} du,$$

et on obtient avec (2.1.7) et (2.1.8) pour tout  $\varepsilon > 0$

$$(2.1.9) \quad \mathbb{P}_f \left\{ \left\| N_h(\hat{\theta}(h) - \theta) \right\|_1 > \varepsilon \right\} \leq \mathbb{P}_f \left\{ \inf_{\|s\|_1 > \varepsilon} \chi_n(s) \leq \chi_n(0) \right\}.$$

Soient  $\nu > 0$  et  $0 < r < \varepsilon/3$  deux paramètres à choisir, après quelques minoration et majorations de  $\chi_n(\cdot)$ , on obtient

$$\begin{aligned} \mathbb{P}_f \left\{ \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0 \right\} &\leq 2\mathbb{P}_f \left\{ \int_{U_n \cap (\|u\|_1 > r)} \|u\|_1 Z_{n,\theta}^{1/m}(u) du > r \nu/4 \right\} \\ &\quad + 2\mathbb{P}_f \left\{ \int_{U_n} Z_{n,\theta}^{1/m}(v) dv < \nu/2 \right\}. \end{aligned}$$

En utilisant l'inégalité de Markov, le premier terme peut être contrôlé par une inégalité exponentielle avec l'hypothèse 3.2. L'hypothèse 3.1 est utilisée pour contrôler le deuxième terme en retranchant  $\nu Z_{n,\theta}^{1/m}(0)$  dans l'inégalité. Ainsi, le contrôle des grandes déviations avec la puissance  $\gamma$  est obtenu pour l'estimateur *bayésien* dans le modèle de *régression générale*. Plus de détails sont donnés dans la preuve 3.5.4 de la proposition 3.

### 2.1.4 Exemples de Modèles avec des Vitesses Différentes

Dans cette partie, nous précisons les modèles de régression (gaussienne, inhomogène de Poisson,  $\alpha$  et multiplicative uniforme) où les hypothèses 3 sont vérifiées. Pour chaque modèle,

on construit un *estimateur bayésien local minimax*. A l'aide de la procédure adaptative, nous donnons un *estimateur bayésien adaptatif optimal* pour chaque modèle.

Nous nous intéressons à trouver des modèles dans lesquels les estimateurs linéaires ne sont pas optimaux (au sens minimax et adaptation minimax). Par exemple, dans les régressions  $\alpha$  et multiplicative uniforme, les vitesses minimax sont respectivement  $\varphi_{n,1+2\alpha}(\beta)$  et  $\varphi_{n,1}(\beta)$  et donc bien meilleures que  $\varphi_{n,2}(\beta)$  la vitesse des estimateurs linéaires. Nous pouvons dire que la famille des estimateurs bayésiens indexés par la fenêtre est plus “riche” que celle des estimateurs linéaires.

Nous ne précisons pas encore, à ce stade, comment vérifier les hypothèses 3. Ce travail est présenté dans la section 3.4. Les démonstrations doivent être faites à la “main”, i.e. pour une densité  $g(\cdot, \cdot)$  donnée, on doit rechercher les constantes qui vérifient les hypothèses 3.

**Régression Gaussienne.** La régression gaussienne est très connue en statistiques non-paramétriques. Beaucoup d'estimateurs linéaires ont déjà été développés et sont optimaux. On veut vérifier que notre estimateur bayésien fait aussi bien que les estimateurs linéaires dans ce cadre très classique. Rappelons le modèle additif gaussien

$$Y_i = f(X_i) + \xi_i, \quad \xi_i \sim \mathcal{N}(0, \sigma^2), \quad i = 1, \dots, n,$$

où  $f \in \mathbb{H}(\beta, L, M)$  et  $\sigma > 0$ . Ici la densité des observations s'écrit

$$g(\cdot, f(X_i)) = g_\xi(\cdot - f(X_i)) = (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{(\cdot - f(X_i))^2}{2\sigma^2} \right\}.$$

On peut vérifier les hypothèses 3 avec  $\gamma = 2$ . Ainsi, l'estimateur bayésien local atteint les vitesses minimax  $\varphi_{n,2}(\beta) = n^{-\frac{\beta}{2\beta+d}}$  et adaptative  $\phi_{n,2}(\beta)$  pénalisée (voir Théorème 6) avec un  $\sqrt{\ln n}$  qui est optimale pour la régression gaussienne. Le choix des fenêtres minimax et adaptatives est donné respectivement dans les sections 2.1.1 et 2.1.2.

**Régression Inhomogène de Poisson.** Nous sommes ici dans le cas discret, avec une densité

$$g(k, f(X_i)) = \mathbb{P}_f(Y_i = k) = \frac{[f(X_i)]^k}{k!} \exp \{-f(X_i)\}, \quad k \in \mathbb{N}, \quad f \in \mathbb{H}_d(\beta, L, M, A).$$

Dans ce cadre, on peut dire que  $Y_i \sim \mathcal{P}(f(X_i))$ . Les hypothèses 3 sont vérifiées avec  $\gamma = 2$  (Lemma 2). Ainsi, l'estimateur bayésien local atteint les vitesses minimax  $\varphi_{n,2}(\beta) = n^{-\frac{\beta}{2\beta+d}}$  et adaptative  $\phi_{n,2}(\beta)$  pénalisée avec un  $\sqrt{\ln n}$  (Théorème 7). Les choix des fenêtres minimax et adaptatives sont données respectivement dans les sections 2.1.1 et 2.1.2. Les vitesses de convergence sont les mêmes que celles de la régression gaussienne. Ce qui est naturel par équivalence asymptotique entre ces deux modèles (voir Anscombe [1948]).

**Régression  $\alpha$ .** La régression  $\alpha$  justifie la construction de l'estimateur bayésien. En effet, la vitesse de convergence dépend de  $\alpha$  et s'améliore dès que  $\alpha < 1/2$ . Rappelons le modèle

$$Y_i = f(X_i) + \epsilon_i, \quad i = 1, \dots, n,$$

où  $\epsilon_i$  a pour densité  $g_\alpha(x) = C(\alpha) \exp\{-|x|^\alpha\}$  avec  $0 < \alpha < 1/2$ .  $C(\alpha)$  est une constante choisie pour que  $g_\alpha$  soit bien une densité. Ici le design  $X_i$  est déterministe uniformément réparti sur  $[0, 1]^d$  et  $f \in \mathbb{H}_d(\beta, L, M)$ .

Nous démontrons que l'estimateur bayésien atteint la vitesse  $\varphi_{n,1+2\alpha}(\beta) = n^{-\frac{\beta}{(1+2\alpha)\beta+d}}$  si l'on connaît  $\beta$ . La procédure adaptative, donnée dans la section 2.1.2, permet de construire un estimateur bayésien adaptatif qui atteint la vitesse  $\phi_{n,1+2\alpha}(\beta)$  pénalisée avec un  $(\ln n)^{\frac{1}{(1+2\alpha)}}$ . On remarque que pour  $\alpha < 1/2$ , on a la vitesse de convergence  $\varphi_{n,1+2\alpha}(\beta) < \varphi_{n,2}(\beta)$ . Ainsi, dans ce modèle, l'estimateur bayésien est “meilleur” (converge plus vite) que les estimateurs linéaires.

Has'minskii et Ibragimov [1981] démontrent, dans le cas paramétrique, que pour  $\alpha > 1/2$  la vitesse est  $n^{-1/2}$  et pour  $\alpha = 1/2$  on a la vitesse  $(n \ln n)^{-1/2}$ . Donc  $1/2$  est la frontière avec la “pire” vitesse  $\varphi_{n,2}(\beta)$ . On remarque que si  $\alpha \approx 0$  alors  $\varphi_{n,1+2\alpha}(\beta) \approx \varphi_{n,1}(\beta)$  où  $\varphi_{n,1}(\beta)$  est apparemment la “meilleure” vitesse que l'on peut atteindre en estimation.

**Régression Multiplicative Uniforme.** Le modèle est

$$Y_i = f(X_i) \times U_i, \quad i = 1, \dots, n,$$

où  $f \in \mathbb{H}_d(\beta, L, M, A)$ . Les variables aléatoires  $(U_i)_{i \in 1, \dots, n}$  sont supposées indépendantes et uniformément distribuées sur  $[0, 1]$  et les points du design  $(X_i)_{i \in 1, \dots, n}$  sont déterministes.

L'étude de ce modèle est l'objet du chapitre 4. L'approche bayésienne développée ne peut pas s'appliquer directement. Premièrement, nous avons construit deux estimateurs  $\hat{A}$  et  $\hat{M}$  de  $A$  et de  $M$ , et introduit l'estimateur bayésien construit sur l'ensemble des coefficients aléatoires. Deuxièmement, le processus  $Z_{n,\theta}$  n'est pas toujours défini (les mesures de probabilité ne sont pas absolument continues). Il nous a fallu choisir de façon judicieuse le polynôme d'approximation  $f_\theta$  avec la condition  $f_\theta \geq f$ .

Ainsi nous prouvons que la vitesse minimax est  $\varphi_{n,1}(\beta) = n^{-\frac{\beta}{\beta+d}}$  (Théorème 10) et que notre estimateur bayésien est minimax (voir Théorème 11, Section 4.2). Dans le cadre adaptatif, nous donnons une procédure plus spécifique que celle présentée ci-dessus, qui nous permet de trouver un estimateur  $\phi_{n,1}(\beta)$ -adaptatif (Théorème 14) où la pénalité est en  $\ln n$ . Nous démontrons à l'aide du critère de Klutchnikoff [2005] que cette famille de normalisations est adaptative optimale (voir Théorème 13). On rappelle qu'il y a toujours un prix à payer pour l'adaptation en estimation ponctuelle.

On rappelle que pour utiliser l'estimateur bayésien, les hypothèses 3 doivent être vérifiées, en particulier pour les exemples ci-dessus. On peut trouver plus de détails dans la section 3.4 (Chapitre 3).

## 2.2 Critère de Huber

Cette partie présente un estimateur “robuste” (peu sensibles aux valeurs extrêmes, en anglais *outliers*) de type M-estimateur que nous appelons estimateur de *Huber*. Comme nous l’avons vu dans la section précédente, l’estimateur bayésien recherche la vitesse optimale pour un modèle donné. L’inconvénient de l’approche bayésienne est qu’il faut supposer connue la densité des observations  $g(\cdot, \cdot)$  pour vérifier les hypothèses 3. Nous proposons un estimateur qui ne dépend pas de la densité  $g_\xi$  des observations (avec  $g_\xi$  vérifiant les hypothèses 1). En revanche, l’estimateur de Huber n’est pas toujours optimal (par exemple régression  $\alpha$  ou régression additive uniforme).

Rappelons dans quel cadre on travaille.

1. Le modèle est celui de la régression additive avec design aléatoire (indépendant du bruit) défini dans la section 1.2.3 avec une densité  $g_\xi(\cdot)$  vérifiant les hypothèses 1,

$$Y_i = f(X_i) + \xi_i, \quad i = 1, \dots, n,$$

2.  $f \in \mathbb{H}_d(\beta, L, M)$  (Définition 1),
3. on utilise l’estimateur de Huber  $\check{f}$  défini dans (1.3.10),
4. on ne regarde pas le point de vue minimax dans ce paragraphe, mais plutôt la vitesse du risque maximal (Définition 6) de notre estimateur de Huber.

L’adaptation dans la régression additive a déjà fait l’objet d’un travail. En effet, [Brown, Cai, et Zhou \[2008\]](#) utilisent la normalité asymptotique de la médiane pour approximer ce modèle par le modèle gaussien avec une méthode de médiane par blocs. Une étape intermédiaire est de projeter les nouvelles observations dans une base d’ondelettes. Ensuite, la méthode de Stein par blocs (voir [Stein \[1981\]](#) et [Cai \[1999\]](#)) est utilisée pour l’adaptation. En revanche, cette approche nécessite des hypothèses plus fortes que les hypothèses 1 qui sont suffisantes pour l’utilisation de l’estimateur de Huber.

[Tsybakov \[1982a, 1982b, 1983, 1986\]](#) fût le premier à introduire l’idée de [Huber \[1981\]](#) pour l’estimation non-paramétrique locale et plus généralement pour l’estimation robuste.

### 2.2.1 Adaptation

Grâce au critère de Huber, on démontre que notre estimateur  $\check{f}$  est  $\phi_{n,2}(\beta)$ -adaptatif sur la collection des espaces de Hölder isotropes  $\left\{ \mathbb{H}_d(\beta, L, M) \right\}_{\beta, L}$  pour l’estimation ponctuelle.

Pour l’estimateur de Huber, citons les papiers de [Hall et Jones \[1990\]](#) et [Härdle et Tsybakov \[1992\]](#), avec les méthodes adaptatives de *validation croisée* et de *Plug-in*. [Reiss, Rozenholc, et Cuenod \[2009\]](#) développe un estimateur adaptatif pour les fonction localement constantes. Dans la suite, nous proposons une généralisation de cette approche à toutes fonctions qui peuvent être, localement, approchées par un objet paramétrique.

Nous supposons que  $\beta \in ]0, b]$  ( $b$  est la borne supérieure de la régularité) et aucune restriction sur les constantes  $L$  et  $M$ . Soit  $\Phi$  la famille de normalisations suivantes :

$$\phi_n(\beta) = \left( \frac{\rho_n(\beta)}{n} \right)^{\frac{\beta}{2\beta+d}}, \quad \rho_n(\beta) = \left( 1 + \frac{2(b-\beta)}{(2\beta+d)(2b+d)} \ln n \right)^{1/2}, \quad \beta \in (0, b].$$

On remarque que  $\phi_n(b) = \varphi_n(b)$  et  $\rho_n(\beta) \sim ((b-\beta)\ln n)^{1/2}$  pour tout  $\beta \neq b$ . Il est facile de démontrer que cette famille de normalisations est adaptative optimale sous le critère de Klutchnikoff [2005] pour la régression gaussienne.

**Construction de l'estimateur  $\Phi$ -adaptatif.** Premièrement, nous déterminons la famille des *estimateurs de Huber locaux*. Ensuite, basée sur la *méthode de Lepski*, nous proposons une règle de sélection dans cette famille à partir des données.

On prend  $\check{f}^h$  l'estimateur défini par (1.3.4), (1.3.9) et (1.3.10). Ainsi, on définit la famille des estimateurs de Huber locaux  $\check{\mathcal{F}}$ . Posons  $h_{\max} = n^{-\frac{1}{2b+d}}$  et

$$h_k = 2^{-k} h_{\max}, \quad k = 0, \dots, \mathbf{k}_n,$$

où  $\mathbf{k}_n$  est le plus grand entier tel que  $h_{\mathbf{k}_n} \geq h_{\min} = n^{-1/d} \ln^{\frac{b}{d(2b+d)}} n$ . Soit

$$(2.2.1) \quad \check{\mathcal{F}} = \left\{ \check{f}^{(k)}(y) = \check{\theta}_{0,\dots,0}(h_k), \quad k = 0, \dots, \mathbf{k}_n \right\}.$$

On pose  $\check{f}^*(y) = \check{f}^{(\hat{k})}(y)$ , où  $\check{f}^{(\hat{k})}(y)$  est sélectionné à partir de  $\check{\mathcal{F}}$  avec la règle de sélection suivante :

$$(2.2.2) \quad \hat{k} = \inf \left\{ k = \overline{0, \mathbf{k}_n} : |\check{f}^{(k)}(y) - \check{f}^{(l)}(y)| \leq C S_n(l), \quad l = \overline{k+1, \mathbf{k}_n} \right\}.$$

Avec les notations suivantes :

$$(2.2.3) \quad C = 2 D_b \sqrt{48 \lambda r d}, \quad S_n(l) = \left[ \frac{1 + l \ln 2}{n(h_l)^d} \right]^{1/2}, \quad l = 0, 1, \dots, \mathbf{k}_n,$$

où  $\lambda > 0$  est la plus petite des valeurs propres de la matrice du design définie dans le lemme 13 et  $D_b$  est le nombre de dérivées partielles d'ordre plus petit ou égal à  $b$  (1.3.2).

**Théorème 3.** Soit  $b > 0$ . Alors, pour toute densité  $g_\xi(\cdot)$  vérifiant les hypothèses 1, pour tout  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$  et  $r \geq 1$

$$\limsup_{n \rightarrow \infty} \phi_{n,2}^{-r}(\beta) R_{n,r} \left[ \check{f}^*(y), \mathbb{H}_d(\beta, L, M) \right] < \infty.$$

**Remarque 6.** Ce théorème correspond au théorème 5 dans le chapitre 3. La constante  $M$  sert à la construction de l'ensemble des coefficients  $\Theta(M)$  (1.3.10), on peut remplacer  $M$  par  $\ln n$  ou l'estimer avec une méthode de médiane locale. La constante  $\lambda$  dépend de  $\int_{-1}^1 g_\xi(z)dz$  (voir Preuve du lemme 13), inconnue en pratique, mais on peut l'estimer à l'aide de l'estimateur

$$\hat{G} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{[-1,1]}(Y_i - \tilde{f}(X_i)),$$

où  $\tilde{f}$  est un estimateur de  $f$  (par exemple l'estimateur de la médiane locale). Il est facile de voir par construction que  $\hat{G}$  estime la quantité  $\mathbb{P}_f(-1 \leq \xi \leq 1) = \int_{-1}^1 g_\xi(z)dz$ . Une solution, très utilisée en pratique, est de calibrer la constante  $C$  avec un jeu de données.

Ce résultat est valable dans le modèle de régression additive sous les hypothèses 1, en particulier l'estimateur de Huber adaptatif  $\check{f}^*(y)$  atteint la vitesse  $\phi_{n,2}(\beta)$  pour les régressions gaussienne, de Cauchy, régression  $\alpha$  et additive uniforme.

Comme pour l'approche bayésienne, nous donnons le résultat sous forme asymptotique. Nous pouvons une nouvelle fois préciser que la borne supérieure peut s'écrire comme la vitesse plus un reste (voir Section 5.4).

Maintenant, nous allons nous intéresser à contrôler les grandes déviations de l'estimateur de Huber.

### 2.2.2 Grandes Déviations

Nous expliquons ici comment obtenir le contrôle des grandes déviations avec l'estimateur de Huber. Les idées de la preuve reposent sur la différentielle du critère de Huber en  $t$ , du processus limite de celle-ci ( $\mathbf{M}$ -estimation) et d'un argument de chaînage combiné avec l'inégalité de Bernstein (Boucheron, Bousquet, et Lugosi [2004]). Tout d'abord, énonçons le résultat.

Pour tout  $n \in \mathbb{N}^*$ ,  $h \in \mathcal{H}_n$ ,  $f \in \mathbb{H}_d(\beta, L, M)$ , alors  $\forall \varepsilon \geq D_b \varepsilon_0(h)/\sqrt{\lambda}$  on a

$$\mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon \right) \leq \Sigma \exp \left\{ - \frac{\left( \varepsilon \sqrt{\lambda} D_b^{-1} - \varepsilon_0(h) \right)^2}{8 + \frac{8}{3\sqrt{nh^d}} \varepsilon \sqrt{\lambda} D_b^{-1}} \right\},$$

où  $\varepsilon_0(h) = (1 \vee b_h N_h)$ . Notons que ce résultat est conditionné par un événement dont son complémentaire a une probabilité exponentiellement petite de se réaliser (pour plus de précisions, voir Section 5.4).

**Idée de la preuve.** Posons la dérivée de la fonction de Huber (1.3.8)

$$(2.2.4) \quad q(z) := Q'(z) = \mathbb{I}_{1,+\infty}[z] - \mathbb{I}_{]-\infty,-1[}(z) + z \mathbb{I}_{[-1,1]}(z), \quad z \in \mathbb{R}.$$



Il est facile de voir que cette fonction est continue et bornée par 1. Notons  $\tilde{D}_h(\cdot)$  la différentielle du critère de Huber  $\tilde{m}_h(\cdot)$  défini par (1.3.9). Remarquons que chaque coordonnée de  $\tilde{D}_h(\cdot)$  est une somme de variables aléatoires bornées et  $\tilde{D}_h(\check{\theta}(h)) = (0, \dots, 0)$  ( $q'$  est continue). Soit  $D_h(\cdot) = \mathbb{E}_f \tilde{D}_h(\cdot)$ , l'idée principale de cette preuve repose sur le fait que  $\tilde{D}_h(t) \xrightarrow[n \rightarrow \infty]{\Delta_h + N_h^{-1}} D_h(t)$  en probabilité uniformément en  $t$ , alors on a  $\check{\theta}(h) \xrightarrow[n \rightarrow \infty]{\Delta_h + N_h^{-1}} \theta$  où  $\theta \in \Theta(M)$  est l'unique solution de  $D_h(\theta) = 0$  et  $\Delta_h$  est le terme de biais est vaut environ  $Lh\beta$ . C'est l'idée générale de la m-estimation : regarder si le processus limite admet comme solution le paramètre que l'on cherche à estimer (voir par exemple le livre de Van de Geer [2000]). Pour obtenir la bonne vitesse de convergence, il suffit d'utiliser les schémas classiques du compromis du biais  $Lh\beta$  et de la variance  $1/\text{sqtr}nh^d$ . On donne le diagramme suivant qui résume ces explications (Figure 2.1).

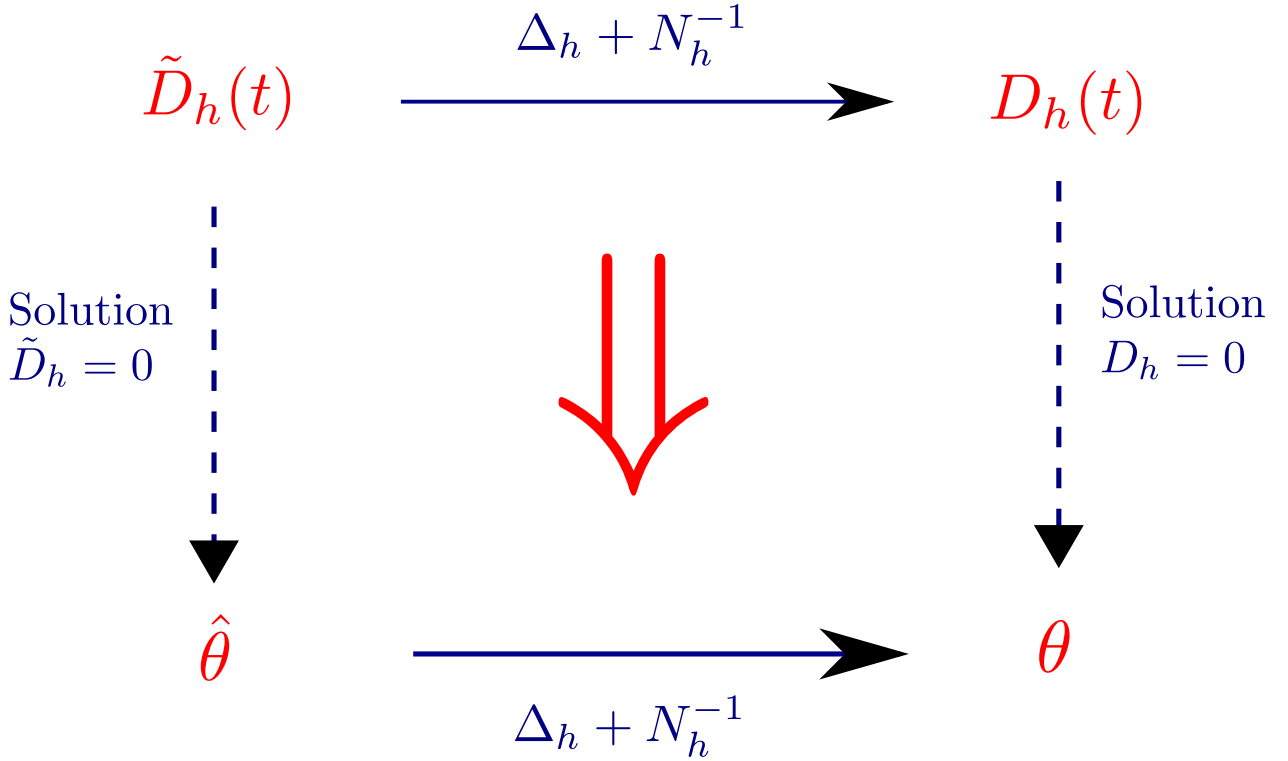


FIGURE 2.1 – Illustration des principe des M-estimateurs

La définition de  $\check{\theta}(h)$  et  $\theta = \theta(f, y, h)$  les coefficients du développement de Taylor impliquent que  $\forall \varepsilon \geq D_b \varepsilon_0(h)/\sqrt{\lambda}$

$$\begin{aligned} \mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon \right) &\leq \mathbb{P}_f \left( N_h |\check{\theta}_0(h) - \theta_0| \geq \varepsilon \right) \\ &\leq \mathbb{P}_f \left( N_h \sqrt{D_b} \|\check{\theta}(h) - \theta\|_2 \geq \varepsilon \right). \end{aligned}$$

En démontrant que la jacobienne de  $D_h(\cdot)$  est inversible, on peut écrire (Lemme 13) que

$$\mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon \right) \leq \mathbb{P}_f \left( N_h \frac{\sqrt{D_b}}{\sqrt{\lambda}} \left\| D_h(\check{\theta}(h)) - D_h(\theta) \right\|_2 \geq \varepsilon \right).$$

Maintenant, rappelons que  $\tilde{D}_h(\check{\theta}(h)) = (0, \dots, 0)$  et  $D_h(\theta) = (0, \dots, 0)$  et

$$\left\| D_h(\check{\theta}(h)) - D_h(\theta) \right\|_2 = \left\| D_h(\check{\theta}(h)) \right\|_2 = \left\| D_h(\check{\theta}(h)) - \tilde{D}_h(\check{\theta}(h)) \right\|_2 \leq \sup_{t \in \Theta(M)} \left\| D_h(t) - \tilde{D}_h(t) \right\|_2.$$

Ainsi

$$\mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon \right) \leq \mathbb{P}_f \left( N_h \frac{\sqrt{D_b}}{\sqrt{\lambda}} \sup_{t \in \Theta(M)} \left\| D_h(t) - \tilde{D}_h(t) \right\|_2 \geq \varepsilon \right)$$

En utilisant un argument de chaînage et l'inégalité de Bernstein (voir Lemme 15), on obtient le résultat

$$(2.2.5) \quad \mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon \right) \leq \Sigma \exp \left\{ - \frac{\left( \varepsilon \sqrt{\lambda} D_b^{-1} - \Delta_h \right)^2}{8 + \frac{8}{3\sqrt{nh^d}} \varepsilon \sqrt{\lambda} D_b^{-1}} \right\}.$$

Pour plus de détails, regarder la proposition 7 dans la section 5.4. On peut voir  $D_h(t) - \tilde{D}_h(t)$  comme une somme de variable aléatoires indépendantes et bornées  $\sum_{i=1}^n \mathcal{W}_t(X_i, \xi_i)$  (processus empirique). Deux nombreux outils ont été développés ces quinze dernières années et peuvent remplacer l'utilisation d'un argument de chaînage et d'une inégalité de concentration classique.

### 2.2.3 Inégalités Maximales pour les processus empiriques

Cette sous-section a pour but de présenter la notion de *processus empirique* et les inégalités de concentration du maxima de ces processus que l'on peut trouver dans la littérature. Ces inégalités sont très utilisées en statistique mathématique notamment pour l'adaptation (Voir l'exemple ci-dessus). Elles sont couramment appelées *inégalités maximales* (en anglais : *maximal inequalities*).

Posons les notations pour définir un processus empirique. Soit  $\mathcal{G}$  une classe de fonctions bornées de  $\Omega \rightarrow \mathbb{R}$  où  $(\Omega, \mathcal{A}, \mathbb{P})$  est un espace de probabilité et  $\mathbb{E}$  désigne l'espérance mathématique par rapport à  $\mathbb{P}$ . Notons  $Z_1, \dots, Z_n$  une suite de variables aléatoires indépendantes à valeurs dans  $\Omega$ . On appelle *processus empirique*  $\sum_{i=1}^n g(Z_i)$  et maxima d'un processus empirique  $Z = \sup_{g \in \mathcal{G}} \sum_{i=1}^n |g(Z_i) - \mathbb{E}g(Z_i)|$ .

Inégalité de Talagrand :

Pour la première fois, Talagrand [1994, 1995, 1996a, 1996b] donne des inégalités de concentration pour le maxima d'un processus empirique. Soient  $b > 0$  tel que  $\|g\|_\infty \leq b$  pour tout  $g \in \mathcal{G}$  et  $v = \mathbb{E} \sup_{g \in \mathcal{G}} \sum_{i=1}^n g^2(Z_i)$ , alors pour tout réel positif

$$\mathbb{P}(Z \geq \mathbb{E}[Z] + x) \leq K \exp \left\{ -\frac{x^2}{2(c_1 v + c_2 b x)} \right\},$$

où  $c_1, c_2$  et  $K$  sont des constantes universelles.

On peut remarquer que si  $\mathcal{G}$  est restreint à une seule fonction, alors on obtient *l'inégalité de Bernstein*

$$\mathbb{P} \left( \sum_{i=1}^n |g(Z_i) - \mathbb{E}g(Z_i)| \geq x \right) \leq 2 \exp \left\{ -\frac{x^2}{2(v + bx/3)} \right\}.$$

Dans ce cas, les constantes correspondantes  $c_1$  et  $c_2$  sont connues explicitement. L'inégalité de Talagrand permet de prouver l'existence d'inégalités exponentielles, mais pose le problème de la connaissance des constantes. En effet, les procédures adaptatives contiennent des paramètres qui, une fois optimisés, dépendent des constantes  $c_1$  et  $c_2$  (Voir le paramètre  $C$  dans la procédure de Lepski). De plus, il n'est pas toujours évident de calculer l'espérance de  $\mathbb{E}[Z]$  (pour plus de détails, voir Massart [2007], Chapitre 6). En fait, cela demande l'utilisation d'un argument de chaînage et donc du calcul de l'entropie de l'espace  $\mathcal{G}$  considéré.

#### Inégalité de Massart 1 :

Pour résoudre les problèmes que pose l'inégalité de Talagrand, Ledoux [1997] propose une démonstration de cette inégalité, en utilisant des inégalités d'entropie sur des mesures produit. Puis Massart [2000] donne une inégalité de type Bennett avec des constantes connues. On peut toujours écrire une inégalité de type Bennett comme une inégalité de type Bernstein qui est souvent privilégiée pour l'adaptation. Soit  $\sigma^2 = \sup_{g \in \mathcal{G}} \sum_{i=1}^n \text{Var}(g(Z_i))$ , alors pour tous réels strictement positifs  $z$  et  $\varepsilon$

$$\mathbb{P}(Z \geq (1 + \varepsilon)\mathbb{E}[Z] + x) \leq \exp \left\{ -\frac{x^2}{2(\kappa\sigma^2 + \kappa(\varepsilon)bx)} \right\},$$

où  $\kappa = 4$  et  $\kappa(\varepsilon) = 2.5 + 32\varepsilon^{-1}$ .

La connaissance des constantes est très importante pour l'adaptation. Le problème, qui vient ensuite, est l'optimisation de ces constantes (les plus petites possibles).

#### Inégalité de Bousquet :

Contrairement à l'inégalité de Massart 1, Bousquet [2002] montre que l'on peut obtenir les constantes optimales (comme dans l'inégalité de Bernstein) avec des variables aléatoires  $Z_1, \dots, Z_n$  indépendantes et identiquement distribuées. Supposons que  $\sup_{g \in \mathcal{G}} \|g\|_\infty \leq 1$ , alors pour tout  $x \geq 0$

$$\mathbb{P}(Z \geq \mathbb{E}[Z] + x) \leq \exp \left\{ -\frac{x^2}{2(\sigma^2 + x/3)} \right\}.$$

Même avec les constantes optimales, il reste un problème de taille à résoudre : le calcul ou la majoration de  $\mathbb{E}[Z]$  (l'espérance du maxima de la somme de variables aléatoires indépendantes et bornées). On peut utiliser un résultat disponible dans le livre de [Van der Vaart et Wellner \[1996\]](#) pour les processus  $(\mathcal{X}_t)_{t \in T}$  dit sub-gaussiens quand ils vérifient

$$\mathbb{P}(|\mathcal{X}_s - \mathcal{X}_t| > x) \leq 2 \exp \left\{ -x^2 / 2d^2(s, t) \right\}, \quad s, t \in T, \quad x > 0,$$

où  $d(\cdot, \cdot)$  est une semi-métrique sur  $T$ . Dans le chapitre 2 (*Empirical processes*), Section 2 (*Maximal inequalities*), corollaire 2.2.8, on y trouve une majoration de l'espérance du maxima du processus  $\mathcal{X}_t$  : pour  $t_0 \in T$  fixé arbitrairement, on a

$$\mathbb{E} \sup_{t \in T} |\mathcal{X}_t| \leq \mathbb{E} |\mathcal{X}_{t_0}| + K \int_0^\infty \sqrt{\ln N(r/2, d)} dr,$$

pour  $K$  une constante universelle et  $N(r, d)$  est le nombre minimal de boules de rayon  $r$  (avec la semi-métrique  $d$ ) nécessaire pour recouvrir l'espace  $T$ .  $\ln N(r/2, d)$  est appelée *entropie* de l'espace  $T$ .

La démonstration de ce corollaire nécessite l'utilisation d'un argument de chaînage et donc le calcul de l'entropie de l'espace considéré. Mais ici, la constante  $K$  est inconnue. Ceci pose encore le problème de l'adaptation dans lequel les paramètres des procédures adaptatives doivent être connues. Dans ce sens, [Massart \[2007\]](#) développa des inégalités maximales qui nécessitent seulement le calcul de l'entropie.

### Inégalité de Massart 2 :

On peut trouver dans le livre de [Massart \[2007\]](#) au Chapitre 6 (*Maximal Inequalities*), Section 6.2 (*Function-indexed empirical processes*), des inégalités exponentielles pour le maxima de processus empiriques (somme de variables aléatoires indépendantes et bornées). Nous expliquerons dans la suite que si le supremum du processus empirique  $\sum_i \mathcal{W}_t(X_i, \xi_i)$  est indexé en  $t$  alors il est nécessaire que  $\mathcal{W}_t(\cdot, \cdot)$  soit continue en  $t$ . Ainsi, le résultat énoncé par [Massart \[2007\]](#) dans le corollaire 6.9 se présente comme suit. L'auteur introduit la fonction  $H(\delta)$  comme l'entropie de  $\mathcal{G}$  avec une notion de *Bracketing* (un peu plus raffinée que l'entropie habituelle). Soit  $\sigma^2 = \sup_{g \in \mathcal{G}} \sum_{i=1}^n \mathbb{E} g^2(Z_i)$  et  $b = \sup_{g \in \mathcal{G}} \|g\|_\infty$ , alors, pour tout  $\varepsilon \in ]0, 1]$  et  $x \geq 0$

$$\mathbb{P}(Z \geq E + x) \leq \exp \left\{ -\frac{x^2}{2(1 + 6\varepsilon)^2 n \sigma^2 + 4bx/3} \right\},$$

où pour une constante absolue  $\kappa = 27$

$$E = \frac{\kappa}{\varepsilon} \sqrt{n} \int_0^{\varepsilon \sigma} \sqrt{H(u) \vee n} du + 2(b/3 + \sigma)H(\sigma)$$

Ici tout le travail a été fait, en particulier la majoration de  $\mathbb{E}Z$  par  $E$  qui contient l'intégrale de l'entropie avec les constantes connues.

Revenons au cas de l'estimateur de Huber, avec le maxima du processus empirique  $\sup_{t \in \Theta} \sum_i \mathcal{W}_t(X_i, \xi_i)$  ( $\Theta = \Theta(M)$  est un compact de  $\mathbb{R}^{D_b}$ ) où  $\mathcal{W}_t(x, z) = q(z + f(x) - f_t(x)) \left(\frac{x-y}{h}\right) \mathbb{I}_{x \in V_h(y)}$ , et  $q$  est la dérivée de la fonction de Huber définie dans (1.3.8). Comme nous l'avons précisé ci-dessus, le calcul de l'entropie sur l'espace  $\mathcal{G} = \mathbb{W}_\Theta = \{\mathcal{W} = \mathcal{W}_t : t \in \Theta\}$  est assez difficile voir impossible. Ainsi, il est souvent choisi de se ramener à l'entropie de l'espace  $\Theta$  qui est plus facile à calculer. Ceci nécessite la continuité de  $\mathcal{W}_t$  en  $t$  et donc de la continuité de la fonction  $q$  (dérivée de Huber) dans notre exemple à nous.

L'inégalité de Massart 2 convient parfaitement pour le contrôle des grandes déviations de l'estimateur de Huber. Mais dans notre démonstration (Chapitre 5, Preuve de la proposition 7), nous avons choisi de faire les calculs à la *main* avec l'utilisation d'un argument de chaînage et l'inégalité de Bernstein (l'inégalité de Hoeffding est suffisante si le design est déterministe).

## 2.3 Expériences numériques

Dans cette section, nous donnons quelques résultats pratiques (avec le logiciel *Matlab*) pour toutes les méthodes développées. Nos simulations sont effectuées sur 3 fonctions  $f_1, f_2$  et  $f_3$ .

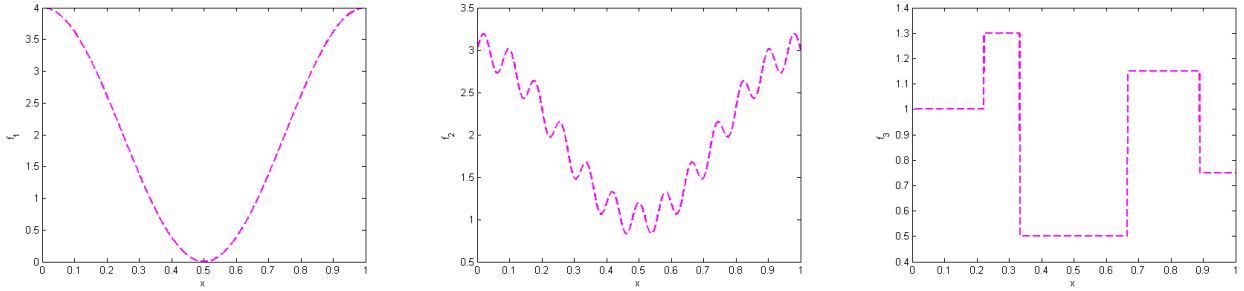


FIGURE 2.2 – Fonctions utilisées pour nos simulations.

Notre étude se limite aux fonctions unidimensionnelles. En fait, les comportements, que l'on peut observer sur les fonctions isotropes à plusieurs variables, sont semblables à ceux observés sur les trois fonctions ci-dessus.

Rappelons que notre estimateur est un polynôme dont on estime les coefficients. Pour des raisons de complexité, nous nous limiterons aux polynômes de degré 1, i.e. que nous n'étudierons que les fonctions  $f \in \mathbb{H}_d(\beta, L, M)$  avec  $\beta \leq 2$ .

Nous testons nos estimateurs dans les différents modèles étudiés dans cette thèse. Pour chacun d'entre eux, nous allons comparer les estimateurs bayésiens et de Huber avec leurs oracles.

**Définition 10.** *L'estimateur  $\hat{f}^*(y) = \hat{f}^{h^*}(y)$  est appelé oracle de la famille  $\{\hat{f}^h(y)\}_{h>0}$  si*

$$h^* = \arg \min_{h>0} \mathbb{E}_f |\hat{f}^h(y) - f(y)|.$$

$h^*$  est appelée *fenêtre oracle*.

Notons que pour chaque point  $y$ , on y associe un oracle. L'oracle peut être interprété comme le meilleur estimateur dans la famille connaissant la fonction  $f$ .

Nous présentons les résultats de la façon suivante : nous calculons les rapports

$$r_B^*(y) = \frac{\mathbb{E}_f |\hat{f}_B^*(y) - f(y)|}{\mathbb{E}_f |\hat{f}_B^{h_B^+}(y) - f(y)|}, \quad r_H^*(y) = \frac{\mathbb{E}_f |\check{f}_H^*(y) - f(y)|}{\mathbb{E}_f |\check{f}_H^{h_H^+}(y) - f(y)|},$$

où  $\hat{f}_B^*(y)$  et  $\check{f}_H^*(y)$  sont respectivement les oracles des familles d'estimateurs bayésiens  $\{\hat{f}^h(y)\}_{h>0}$  et de Huber  $\{\check{f}^h(y)\}_{h>0}$ .  $h_B^+$  et  $h_H^+$  sont respectivement les fenêtres adaptatives des estimateurs bayésien et de Huber, données par les procédures (2.1.5) et (2.2.2). L'espérance de la

fonction de perte est approchée par une méthode de Monte-Carlo avec  $T = 10000$  itérations. On construit la suite finie d'éléments  $y_1, y_2, \dots, y_v$  uniformément répartis sur l'intervalle  $[0, 1]$ . On appelle *oracle bayésien* / *estimateur bayésien* et *oracle de Huber* / *estimateur de Huber* les rapports suivants :

$$r_B^* = \frac{1}{v} \sum_{j=1}^v r_B^*(y_j), \quad r_H^* = \frac{1}{v} \sum_{j=1}^v r_H^*(y_j).$$

Nous choisirons  $v = 200$  pour l'implémentation. Ces rapports permettent d'évaluer la qualité des procédures adaptatives que nous avons introduit dans cette thèse (voir Sections 2.1 et 2.2). Rappelons que ces procédures contiennent un paramètre de seuillage, noté  $C$  (voir (2.2.3) et (2.1.6)). Nous proposons un choix théorique connu de ce paramètre, mais en pratique, on peut le choisir de façon optimale. La figure 2.3 illustre le choix optimal de ce paramètre. En effet, il suffit de le calibrer par rapport à un jeu de données ou de fonctions *fantômes* pour calculer l'oracle et regarder la distance qui sépare notre estimateur à cet oracle.

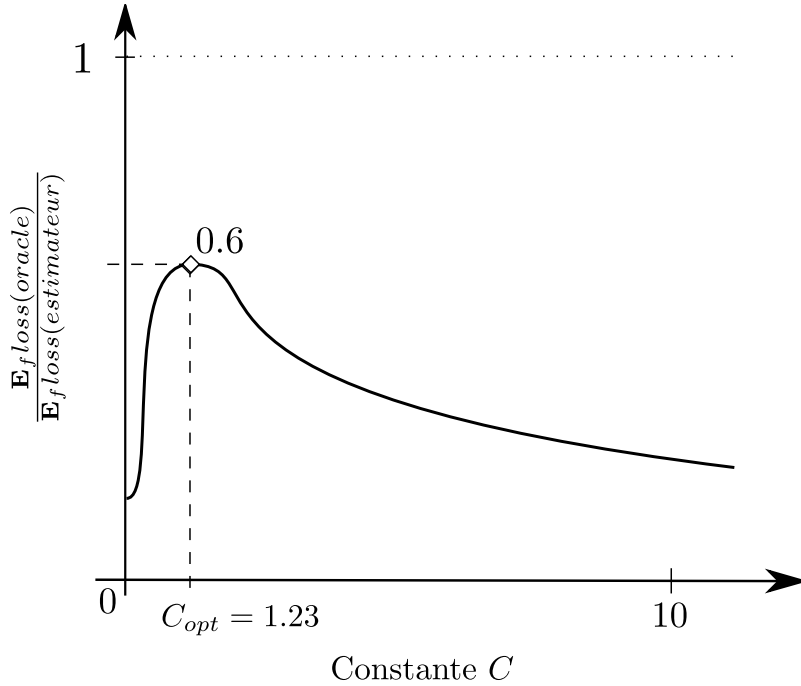


FIGURE 2.3 – Illustration de la calibration de la constante  $C$  à travers un jeu de données connu. Ici, le modèle est celui de la régression gaussienne (avec la fonction  $f_1$ ), avec  $n = 1000$  observations. Nous itérons  $T = 10000$  séries d'observations pour approximer l'espérance du risque. La fenêtre de l'estimateur bayésien adaptatif est choisie avec la procédure donnée dans la section 2.1. Nous avons remarquer que le  $C_{th}$  est cinq à dix fois plus grand que le  $C_{opt}$ .

Nous constatons ici une nette différence entre le risque de notre estimateur et celui de l'oracle (avec un rapport légèrement supérieur à 1/2). Ceci met en évidence les points faibles de l'approche adaptative au sens minimax et de l'estimation ponctuelle (prix à payer). En effet, Nous aimerions que notre rapport  $r^*$  soit le plus proche possible de 1, au minimum en  $C_{opt}$ . Une perspective serait d'utiliser une méthode adaptative pour obtenir des inégalités oracle (Voir Section 2.4).

Nous donnons aussi, quand cela est possible, les rapports *oracle linéaire / estimateur bayésien* et *oracle linéaire / estimateur de Huber*, notés  $r_{LB}^*$  et  $r_{LH}^*$ , où l'oracle Linéaire est celui d'une famille d'estimateurs à noyau classique (noyau d'Epanechnikov). Cela nous permet de comparer la qualité des nouveaux estimateurs, développés dans cette thèse, avec les plus classiques comme les estimateurs linéaires.

Rappelons que la complexité des estimateurs bayésien et de Huber est plus grande que celle des estimateurs linéaires. Pour le critère de Huber qui est convexe, cela peut se faire avec une méthode de descente du gradient. En revanche, le critère bayésien n'est pas du tout évident à minimiser surtout si la densité des observations est discontinue. Ceci est le point faible de ces nouveaux estimateurs (non explicites). En revanche, avec le même d'observations, l'estimateur bayésien est plus rapide dans certains cas, ce qui entraîne un besoin moins important en observations. En pratique, la construction d'observations peut coûter cher en temps et en argent. Un autre aspect est de considérer les modèles où les estimateurs linéaires ne sont pas consistants (par exemple, modèle de Cauchy). Ainsi, nous pouvons utiliser l'estimateur de Huber pour construire un estimateur au moins consistant.

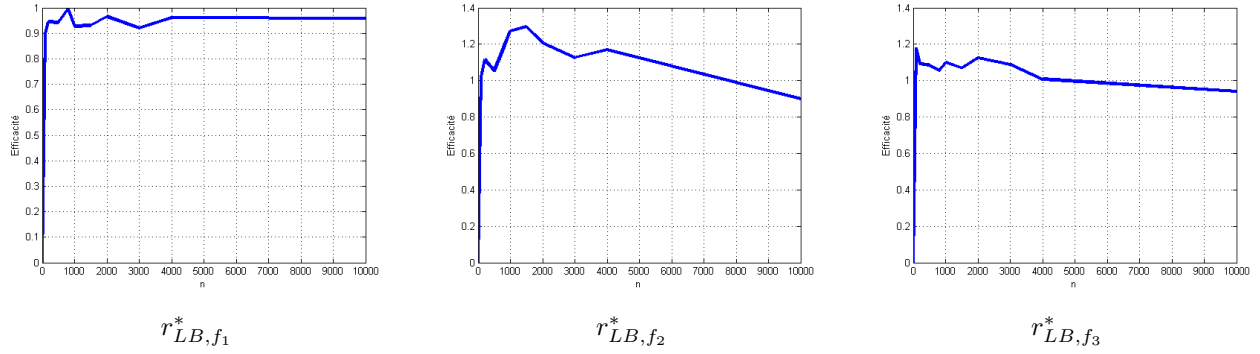
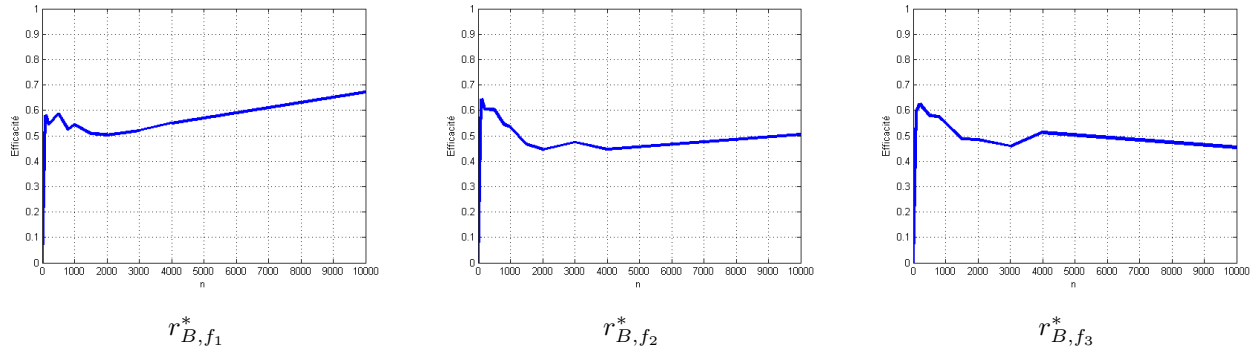
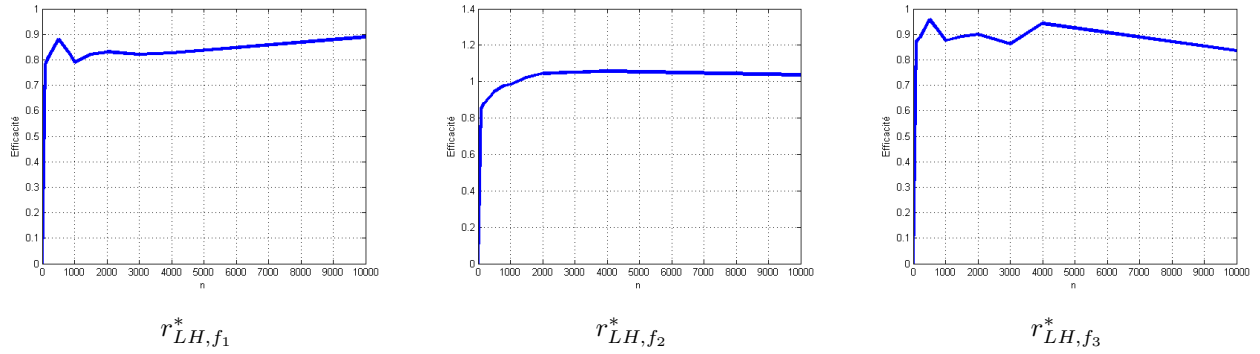
**Légende.** Pour chaque modèle, nous illustrons nos simulations en traçant les courbes  $r_B^*$ ,  $r_H^*$ ,  $r_{LB}^*$  et  $r_{LH}^*$  en fonction de  $n$  (jusqu'à  $n = 10000$ ). Ceci est effectué pour chaque fonction de régression  $f_1$ ,  $f_2$  et  $f_3$ . Dans le même temps, nous calculons la probabilité  $\mathbb{P}_f(h^* < h^+)$  pour chaque estimateur adaptatif. Cette probabilité correspond aux grandes déviations de l'estimateur en question. Le choix du paramètre de seuillage  $C$  est fait de façon optimale (voir Figure 2.3) et est donné pour chaque modèle.

## Régression Gaussienne

On peut voir sur la figure 2.4 que l'oracle bayésien fait aussi "bien" que l'oracle linéaire. Ce qui confirme les résultats théoriques obtenus pour ce modèle (voir Section 2.1.4). A une constante près, les deux estimateurs ont la même vitesse  $\phi_{n,2}$ . La famille d'estimateurs bayésiens est donc aussi "riche" que celle des estimateurs linéaires.

On constate sur cette illustration (Figure 2.5) que le risque de l'estimateur adaptatif est deux fois plus grand que celui de l'oracle. Ceci est classique en estimation ponctuelle du fait du *prix à payer* ( $\ln^{1/\gamma} n$ ) pour l'adaptation (Voir Théorème 2). Il faudrait prendre un nombre d'observations  $n \gg 10000$  (bien plus grand que 10000) pour voir le rapport tendre vers 1. Ce phénomène s'explique en particulier par le fait que les grandes déviations sont environ égales à 4%.



FIGURE 2.4 – Rapport *oracle linéaire* / *oracle bayésien* (modèle gaussien).FIGURE 2.5 – Rapport *oracle bayésien* / *bayésien adaptatif* :  $C_{opt} = 1.2$  and  $\mathbb{P}(\hat{h} < h^*) \simeq 0.0411$  (modèle gaussien).FIGURE 2.6 – Rapport *oracle Linéaire* / *oracle Huber* (modèle gaussien).

De même que l'oracle bayésien, l'oracle de Huber fait aussi bien que l'oracle linéaire (Figure 2.6). Nous signalons que le rapport est légèrement en-dessous de 1 (plutôt 0.9). Ceci est dû aux constants de la borne supérieure pour l'estimateur de Huber. Cette constante est en effet plus grande que celle que l'on peut obtenir pour un estimateur linéaire. Ce

phénomène se retrouve entre la médiane et la moyenne empirique.

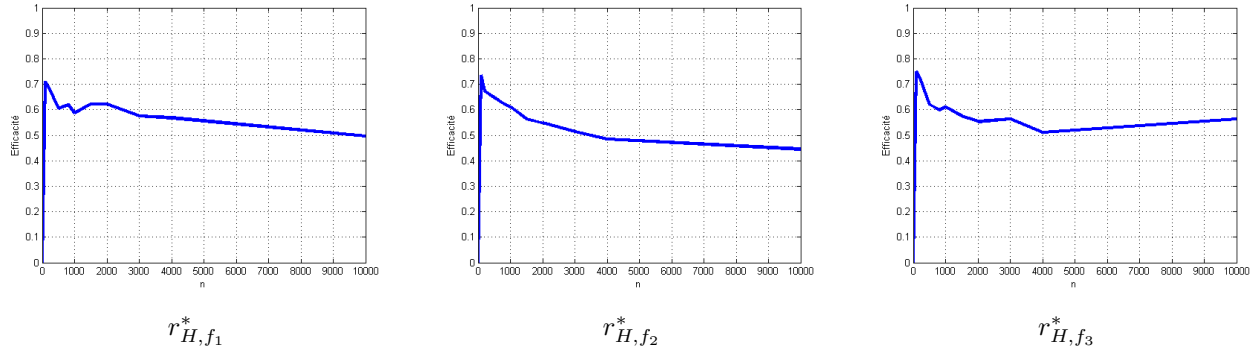


FIGURE 2.7 – Rapport *oracle* Huber / Huber adaptatif :  $C_{opt} = 1.2$  and  $\mathbb{P}(\hat{h} < h^*) \simeq 0.065$  (modèle gaussien).

Nous constatons le même effet que dans le cas précédent. L'estimateur adaptatif de Huber est deux fois moins bon que l'oracle de Huber (Figure 2.7). Comme cette remarque a déjà été faite pour l'estimateur bayésien adaptatif, nous pensons que cela est dû au modèle avec des grandes déviations égales à 6.5%. La forme de l'estimateur en question ne joue aucun rôle, nous mettrons en avant la sensibilité de la méthode de Lepski par rapport au modèle. Nous verrons dans la suite d'autres modèles où les valeurs des rapports *oracle* / *estimateur adaptatif* changent.

### Régression Inhomogène de Poisson

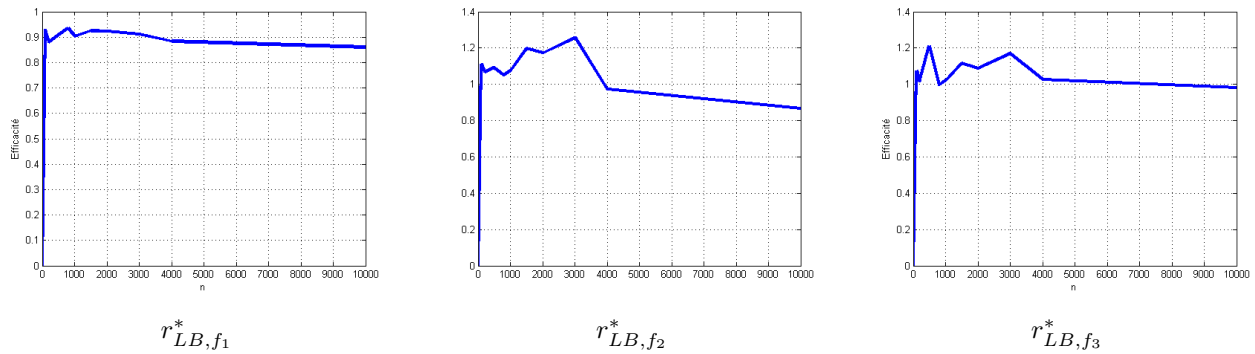


FIGURE 2.8 – Rapport *oracle* linéaire / *oracle* bayésien (Modèle de Poisson).

Le modèle de Poisson est asymptotiquement équivalent au modèle gaussien (voir [Anscombe \[1948\]](#)). De ce fait, on retrouve les mêmes caractéristiques que le modèle gaussien. L'oracle bayésien est aussi bon que l'oracle linéaire (Figure 2.8). Le rapport *oracle* bayésien / bayésien adaptatif vaut environ 1/2 avec  $\mathbb{P}(\hat{h} < h^*) \simeq 0.0480$  (Figure 2.9).

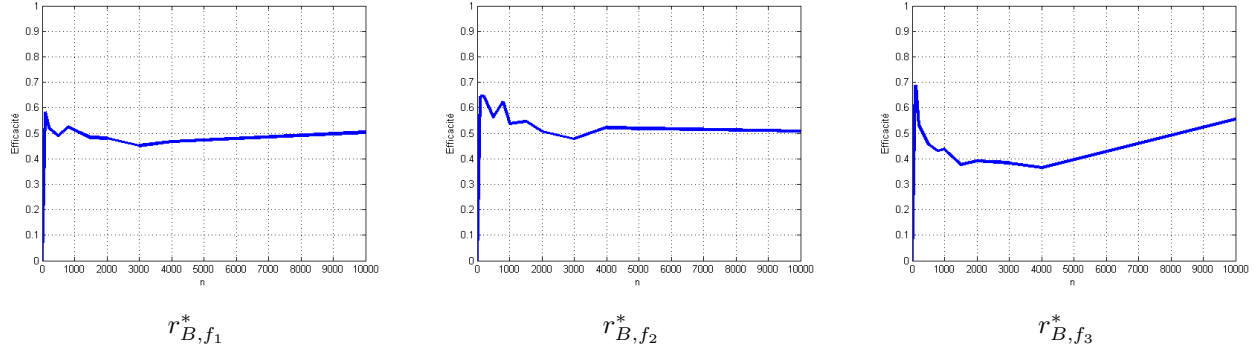


FIGURE 2.9 – Rapport *oracle bayésien* / *bayésien adaptatif* :  $C_{opt} = 1.2$  and  $\mathbb{P}(\hat{h} < h^*) \simeq 0.0480$  (Modèle de Poisson).

### Régression $\alpha = 1/4$

Nous avons choisi pour nos simulations  $\alpha = 1/4$ .

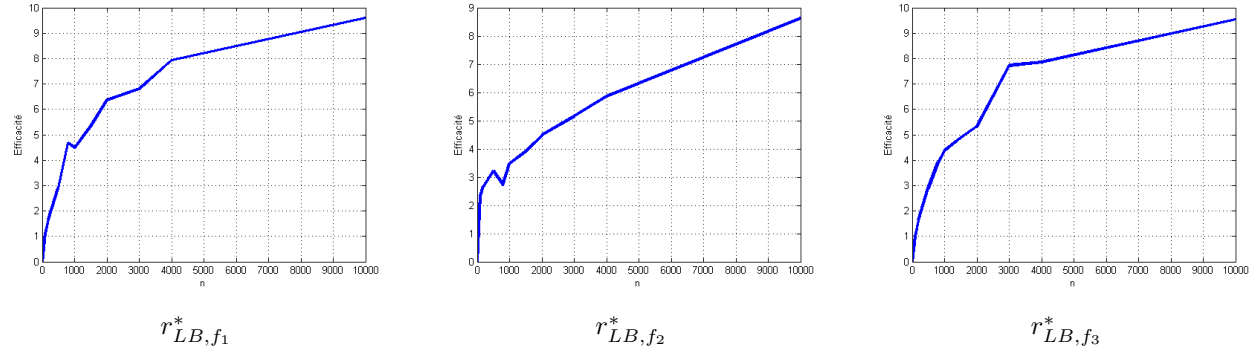


FIGURE 2.10 – Rapport *oracle linéaire* / *oracle bayésien* (modèle  $\alpha$ ).

Comme nous l'avons vu dans la section 2.1.4, la vitesse atteinte par l'estimateur bayésien est  $\phi_{n,1.5}$ . Elle est donc meilleure que la vitesse des estimateurs linéaires  $\phi_{n,2}$ . Nous constatons ce phénomène sur les simulations du rapport *oracle linéaire* / *oracle bayésien* (Figure 2.10). Ici, le rapport peut atteindre la valeur de 10 pour  $n = 10000$ , i.e que l'oracle bayésien est dix fois plus rapide que l'oracle linéaire.

L'étude de la performance de la méthode adaptative (méthode de Lepski) pour la régression  $\alpha$ , révèle que celle-ci est assez faible. En effet, le rapport *oracle bayésien* / *bayésien adaptatif* est quasiment égal à 0.1 (Figure 2.11). Ce qui veut dire que le risque de l'estimateur adaptatif est 10 fois moins bon que celui de l'oracle. On l'explique par le fait que les grandes déviations sont relativement élevées (par rapport aux modèles précédents, égales à 17%). Ce modèle comporte un certain nombre de valeurs extrêmes qui peut nous laisser penser que celles-ci mettent à mal la méthode de Lepski. En effet, nous avons mentionné dans la section 1.5.2 que cette méthode s'arrête trop tôt en présence de valeurs extrêmes.

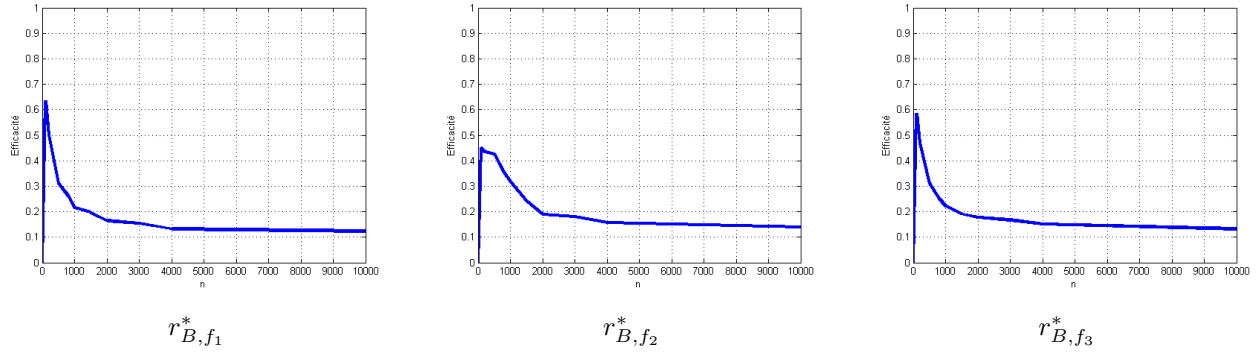


FIGURE 2.11 – Rapport *oracle bayésien* / *bayésien adaptatif* :  $C_{opt} = 2$  and  $\mathbb{P}(\hat{h} < h^*) \simeq 0.1703$  (modèle  $\alpha$ ).

### Régression Multiplicative Uniforme

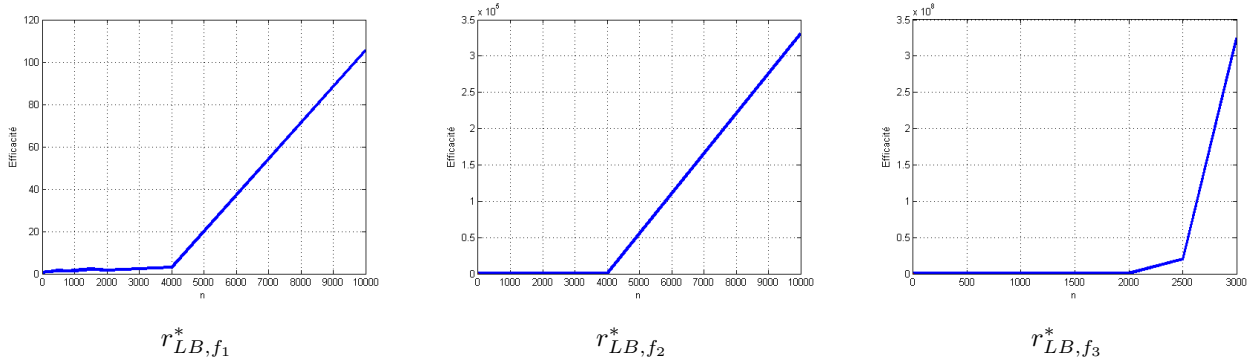


FIGURE 2.12 – Rapport *oracle linéaire* / *oracle bayésien* (modèle multiplicatif uniforme).

Dans cette illustration (Figure 2.12), on constate un net avantage pour l'oracle bayésien. En effet, pour la fonction  $f_1$ , le rapport *oracle linéaire* / *oracle bayésien* atteint 250. Ensuite, pour les fonctions irrégulières  $f_2$  et  $f_3$ , nous avons des valeurs qui sont de l'ordre de  $10^5$  et  $10^8$ . Remarquons que pour la fonction  $f_3$ , nous nous sommes arrêté à  $n = 3000$ , à cause des valeurs trop petites du risque de l'oracle bayésien que le logiciel *Matlab* arrondit à 0. Deux remarques permettent d'expliquer les bonnes performances de l'estimateur bayésien. Premièrement, la vitesse atteinte  $\phi_{n,1}$  est meilleure que  $\phi_{n,2}$ . Deuxièmement, pour les fonctions irrégulières, l'estimateur bayésien est plus performant pour l'estimation au voisinage des points de discontinuité, tandis que l'estimateur linéaire devient très mauvais pour ce genre de fonctions.

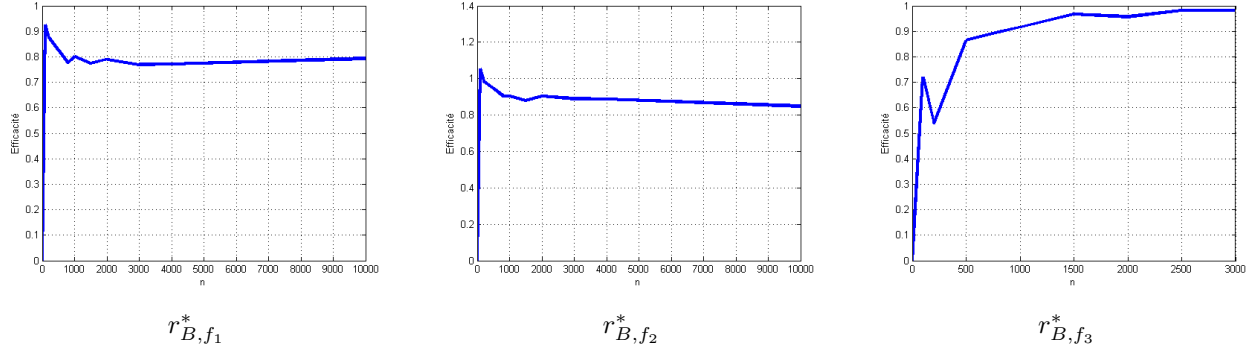


FIGURE 2.13 – Rapport *oracle bayésien / bayésien adaptatif* :  $C_{opt} = 2$  and  $\mathbb{P}(\hat{h} < h^*) \simeq 0$  (modèle multiplicatif uniforme).

En constatant que l'estimateur adaptatif fait aussi bien que son oracle (Figure 2.13), nous pouvons souligner une nouvelle fois la robustesse de notre estimateur, car les grandes déviations sont nulles. Ce dernier modèle souligne une nouvelle fois, l'intérêt de trouver la vitesse optimale et l'estimateur qui l'atteigne. Dans ce sens, l'estimateur bayésien joue pleinement son rôle.

### Régression de Cauchy

Pour le modèle de Cauchy, il n'est pas possible de construire les estimateurs linéaires, du fait qu'ils ne sont pas consistants. Nous donnons, pour évaluer la performance de notre estimateur de Huber, le risque de l'estimateur adaptatif en moyenne pour les trois fonctions ( $n = 10000$ ) :

$$\frac{1}{v} \sum_{j=1}^v \mathbb{E}_f |\check{f}_H^{h^+}(y_j) - f(y_j)| \simeq 10^{-1}.$$

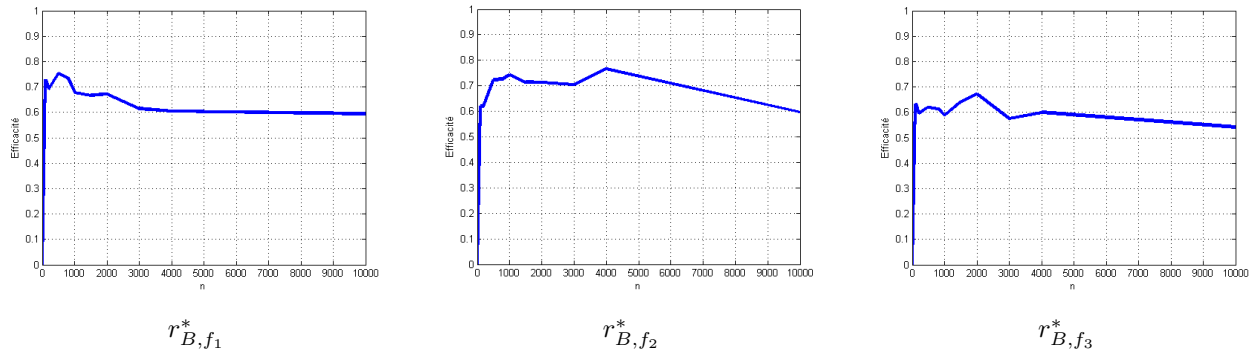


FIGURE 2.14 – Rapport *oracle de Huber / Huber adaptatif* :  $C_{opt} = 3$  and  $\mathbb{P}(\hat{h} < h^*) \simeq 4.7\%$  (modèle de Cauchy).

Ceci est une bonne performance si on prend en compte les constantes de la borne supérieure et la vitesse qui vaut environ  $1/\sqrt{10000} = 10^{-2}$  pour les fonctions discontinues.

Pour la partie adaptation, nous constatons que le rapport *oracle de Huber / Huber adaptatif* vaut 0.6 pour les trois fonctions (Figure 2.14). Nous soulignons le caractère robuste de l'estimateur de Huber qui, malgré un bruit de Cauchy (avec un grand nombre de valeurs extrêmes), rend stable la procédure de Lepski avec seulement  $\mathbb{P}(\hat{h} < h^*) \simeq 0.047\%$ .

## 2.4 Perspectives

Nous avons présenté dans cette thèse deux nouvelles méthodes d'estimation : l'une qui recherche la vitesse optimale (bayésienne) et l'autre qui est robuste et qui ne dépend pas de la densité du bruit (Huber). Nous donnons, dans la suite, une quinzaine de problèmes ouverts relatifs à ces deux approches.

### 2.4.1 Approche Bayésienne

1. L'approche bayésienne est déduite des travaux de [Has'minskii et Ibragimov \[1981\]](#) pour les modèles paramétriques. Nous avons insisté sur le fait que le processus  $Z_{n,\theta}(\cdot)$  (rapport de vraisemblance) doit vérifier les hypothèses 3. Mais en réalité, on peut prendre n'importe quel processus qui vérifie ces hypothèses. Par exemple, dans le modèle de la régression additive uniforme définie par

$$Y_i = f(X_i) + U_i, \quad U_i \sim \mathcal{U}_{[-1,1]}, \quad i = 1, \dots, n.$$

Avec ce modèle, nous ne pouvons pas définir le processus  $Z_{n,\theta}(\cdot)$  car les mesures de probabilité ne sont absolument pas continues.

**Problème Ouvert 1.** *Choisir un processus  $Z'_{n,\theta}$  qui vérifie les hypothèses 3 et qui permet de démontrer que l'estimateur bayésien atteint la vitesse minimax  $n^{-\frac{\beta}{\beta+d}}$  pour la régression additive uniforme ci-dessus. Ensuite, montrer que l'estimateur bayésien adaptatif atteint la vitesse  $\phi_{n,\gamma}$ .*

2. Pour le modèle multiplicatif uniforme, nous avons émis des hypothèses supplémentaires  $0 < A \leq f \leq M$ . La borne supérieure de la fonction  $f$  est une hypothèse classique. En revanche, l'hypothèse de positivité de celle-ci est absolument nécessaire pour l'approche bayésienne mais non justifiée.

**Problème Ouvert 2.** *Pour le modèle de régression multiplicatif uniforme, affaiblir l'hypothèse de positivité de la fonction  $f$  à estimer, i.e. supposer seulement que  $\|f\|_\infty < \infty$ .*

3. Nous pensons que l'on peut atteindre la vitesse  $n^{-\frac{\beta}{\beta+d}}$  pour tout modèle où la densité des observations admet un nombre fini de points de discontinuité. Cela a déjà été démontré dans le cas paramétrique par [Has'minskii et Ibragimov \[1981\]](#) (Chapitre 5, Section 1).

**Problème Ouvert 3.** *Soit le modèle de régression générale (voir 1.2.2), où la densité  $g$  admet un nombre fini de points de discontinuité. Démontrer que*



la vitesse  $n^{-\frac{\beta}{\beta+d}}$  est atteinte par l'estimateur bayésien, i.e. que les hypothèses 3 sont vérifiées avec  $\gamma = 1$ .

4. Dans cette thèse, nous avons supposé une hypothèse d'indépendance des observations, mais celle-ci n'est pas nécessaire. On pourra alors traiter des modèles avec des observations dépendantes.

**Problème Ouvert 4.** *Utiliser l'approche bayésienne développée dans cette thèse pour des modèles où les observations ne sont pas indépendantes (par exemple, faiblement dépendantes).*

5. L'étude des modèles hétéroscédastiques (i.e. la variance n'est pas constante,  $Y_i = f(X_i) + \sigma(X_i)\xi_i$ ) est une perspective intéressante. Nous pensons que l'estimateur bayésien peut être appliqué à ce genre de modèles.

**Problème Ouvert 5.** *Peut-on utiliser l'estimateur bayésien dans des modèles hétéroscédastiques ?*

6. Dans notre travail, nous nous sommes restreints aux espaces isotropes. Cela est dû à l'adaptation et à la méthode de Lepski. Il est possible d'atteindre la vitesse  $n^{-\frac{\bar{\beta}}{\gamma\bar{\beta}+1}}$  ( $\bar{\beta}$  est la moyenne harmonique) sur les espaces de Hölder anisotropes avec l'estimateur bayésien dans le modèle de régression générale. En revanche, pour l'adaptation sur ces espaces, la méthode développée par Kerkycharian, Lepski, et Picard [2001] ne convient pas pour les estimateurs non-linéaires.

**Problème Ouvert 6.** *Existe-t-il une méthode adaptative capable d'estimer de façon optimale (atteindre la vitesse  $\phi_{n,\gamma}(\bar{\beta})$ ) les fonctions anisotropes avec l'estimateur bayésien ?*

7. Le point de vue minimax adaptatif est souvent vu comme pessimiste. La recherche d'inégalités de type oracle permet d'éliminer ce problème.

**Problème Ouvert 7.** *Obtenir des inégalités oracle pour la famille des estimateurs bayésiens indéxés par la fenêtre. Par exemple, l'approche de Golden-shluger et Lepski [2009a] peut convenir pour l'étude des fonctions isotropes.*

8. Les hypothèses 3 sont suffisantes pour montrer que l'estimateur bayésien est adaptatif et optimal (dans certains modèles). La question peut se formuler comme suit.

**Problème Ouvert 8.** *Existe-t-il des hypothèses nécessaires pour montrer que l'estimateur bayésien est adaptatif optimal ? Plus généralement, peut-on trouver les hypothèses nécessaires pour l'existence d'un estimateur adaptatif optimal ?*

9. Nous avons montré que l'estimateur bayésien atteignait la vitesse  $\phi_{n,\gamma}$  pour la régression générale sous les hypothèses 3. De façon générale, il n'est pas prouvé que cette vitesse est optimale. On peut seulement montrer son optimalité modèle par modèle. Une perspective est de le démontrer pour n'importe quel modèle.

**Problème Ouvert 9.** *Démontrer, avec le critère de Klutchnikoff [2005], que la vitesse  $\phi_{n,\gamma}$  est optimale pour tout modèle sous les hypothèses 3.*

10. Un autre problème ouvert est celui des hypothèses. Nous constatons que les conditions sur la densité  $g(.,.)$  se font à travers les hypothèses 3 sur le rapport de vraisemblance  $Z_{n,\theta}$ . Nous aimerions trouver les conditions nécessaires sur la densité (à l'aide de la distance d'Hellinger) pour construire un estimateur bayésien adaptatif optimal.

**Problème Ouvert 10.** *Peut-on donner des hypothèses sur la densité  $g$  qui impliquent les hypothèses 3 ?*

11. L'approche bayésienne consiste à mettre une loi *a priori* sur le paramètre à estimer. Dans ce cas, le critère bayésien devient  $\int \|t - u\| d_a(u) du$ , où  $d_a(u)$  remplace la vraisemblance  $(L_h(u))^{1/m}$  utilisée dans cette thèse. Il est possible de démontrer nos résultats pour n'importe quel *a priori*  $d_a(u)$  sous les hypothèses 3.

**Problème Ouvert 11.** *Remplacer  $Z_{n,\theta}$  par un *a priori* et donner une méthode pour choisir l'*a priori* optimal sous les hypothèses 3 (voir par exemple Ghosal, Lember, et Van der Vaart [2008] et Van der Vaart et Van Zanten [2009]).*

### 2.4.2 Critère de Huber

L'approche de Huber est nouvelle pour l'adaptation. Seuls les résultats de Reiss, Rozenholc, et Cuenod [2009] traitent du sujet pour l'estimateur localement constant et non localement polynômial avec un degré quelconque.

1. Pour cette approche, nous nous sommes aussi restreints aux espaces isotropes. En effet, les approches développées par Kerkycharian, Lepski, et Picard [2001], Goldenshluger et Lepski [2008] et Goldenshluger et Lepski [2009a] ne conviennent pas pour les estimateurs non-linéaires (on ne peut pas comparer les biais avec des estimateurs non-linéaires). Nous avons espoir de trouver de nouvelles méthodes qui permettent de traiter les fonctions anisotropes. La recherche d'inégalités oracle est aussi une perspective.

**Problème Ouvert 12.** *Existe-t-il une méthode adaptative pour obtenir des inégalités oracle pour les fonctions anisotropes avec la famille des estimateurs de Huber indexés par la fenêtre ?*

2. Comme pour l'approche bayésienne, nous pensons pouvoir utiliser l'approche de *Huber* pour les modèles hétéroscedastiques. En effet, le contrôle des grandes déviations repose sur l'inégalité de Bernstein qui peut être établie dans le cas hétéroscedastique. En revanche, l'indépendance des observations est absolument nécessaire pour utiliser cet estimateur.

**Problème Ouvert 13.** *Peut-on utiliser l'estimateur de Huber dans des modèles hétéroscedastiques ?*

3. Nous pensons aussi que les hypothèses 1 peuvent être affaiblies. Notamment, il serait intéressant d'étudier le problème ouvert suivant.

**Problème Ouvert 14.** *Peut-on obtenir les mêmes résultats avec l'approche de Huber avec seulement l'hypothèse de symétrie du bruit ? Dans le cas où le bruit est symétrique et  $g(0) = 0$ , peut-on trouver une vitesse de convergence meilleure que celle obtenue jusqu'ici par l'estimateur de Huber ? (l'information de Fisher est égale à l'infinie).*

4. Il est important de prouver que la vitesse  $\phi_{n,2}$  (atteinte par l'estimateur de Huber), est optimale. Par exemple, pour le bruit gaussien ou de Cauchy, il est possible de le démontrer (avec les travaux de [Klutchnikoff \[2005\]](#), dans le cas minimax voir [Tsybakov \[2008\]](#), Chapitre 2). En revanche, dans le modèle du bruit additif uniforme, cette vitesse n'est pas optimale, car déjà dans le cas minimax, on a  $n^{-\frac{\beta}{\beta+d}}$ .

**Problème Ouvert 15.** *A l'aide des travaux de [Klutchnikoff \[2005\]](#), donner les conditions nécessaires sur  $g_\xi$ , pour que la vitesse adaptative  $\phi_{n,2}$  soit optimale.*

5. Remarquons que la fonction de Huber mélange la norme  $\ell_1$  et la norme  $\ell_2$ . La norme  $\ell_1$  est là pour contrôler les valeurs extrêmes (cela utilise les avantages de la médiane) et la norme  $\ell_2$  est utilisée autour de 0 pour rendre le critère continu dans son comportement. La taille du voisinage autour de 0 où est utilisée la norme  $\ell_2$  peut être changée (dans notre cas, la taille est égale à 1). Cela ne change rien en théorie mais en pratique, ce voisinage est à calibrer. Les estimateurs robustes sont de plus en plus demandés en imagerie pour leur stabilité (voir [Arias-Castro et Donoho \[2009\]](#) ou [Astola, Egiazarian, Foi, et Katkovnik \[2010\]](#)). Développer cette approche pour l'imagerie (ou d'autres domaines) est une perspective intéressante.

**Problème Ouvert 16.** *En pratique, est-ce que la taille du voisinage, de la norme  $\ell_2$  dans le critère de Huber, améliore la performance ?*

6. Au vue des preuves du chapitre 5, il est possible d'utiliser d'autre critère que celui de Huber.

**Problème Ouvert 17.** *Trouver les conditions nécessaires sur le critère (par exemple, convexité, dérivée bornée, etc.) pour construire un estimateur adaptatif qui atteigne la vitesse  $\phi_{n,2}$  dans le modèle de régression additive (voir Section 1.2.3).*



# Chapter 3

## General Locally Bayesian Approach

Ce chapitre traite du modèle de régression générale présenté dans la section 1.2. Nous développons un *estimateur Bayésien* déjà présenté en section 1.3.1. Le contenu de ce chapitre peut être trouvé dans l'article de Chichignoud [2010b], qui généralise le chapitre 4 et l'article de Chichignoud [2010a]. Sous certaines hypothèses 1 sur la densité du modèle, nous montrerons que notre *estimateur Bayésien* est optimal au sens minimax et minimax adaptatif (voir sections 3.2 et 3.3). Dans la section 3.4 nous présentons des exemples de modèles de régression, nous donnons en particulier les preuves qui permettent de vérifier les hypothèses 1 pour ces exemples et ainsi établir les vitesses de convergences minimax pour chaque modèle. Le contenu de ce chapitre généralise les résultats du chapitre 4. Les résultats de ce chapitre permettent de conclure le fait suivant : la famille des *estimateurs Bayésiens* est plus “riche” que les *estimateurs linéaires*.

### 3.1 Introduction

Let statistical experiment

$$\mathfrak{G}_n = \left\{ [0, 1]^d \times \mathbb{R}^n, \mathbb{B}^{(n)}, \mathbb{P}_f^{(n)}, f \in \mathfrak{F} \right\}, \quad n \in \mathbb{N}^*,$$

be generated by independent random observations

$$(3.1.1) \quad \mathcal{Z}_n = (X_i, Y_i)_{i=1, \dots, n},$$

where the  $d$ -dimensional vector  $X_i \in [0, 1]^d$  can be viewed as a location in time or space and  $Y_i$  as the “observation at  $X_i$ ”. Here  $f : [0, 1]^d \rightarrow \mathbb{R}$  is unknown function and we are interested in estimating  $f$  at a given point  $y \in [0, 1]^d$  from observation  $\mathcal{Z}_n$ .

The design points  $(X_i)_{i=1, \dots, n}$  are deterministic and without loss of generality, we will assume that

$$X_i \in \left\{ 1/n^{1/d}, 2/n^{1/d}, \dots, 1 \right\}^d, \quad i = 1, \dots, n.$$

The design can be considered random and the proof is the same. But, we choose the deterministic design to simplify the notation.

Our model assumes that the values  $X_i$  are given and a distribution  $g(\cdot, f_i)$  of each  $Y_i$  is determined by a parameter  $f_i$ , which may depend on the location  $X_i$ ,  $f_i = f(X_i)$ . The estimation problem is to reconstruct  $f(y)$  from the data  $\mathcal{Z}_n$ . Let us illustrate this set-up by few special cases.

1. *Gaussian regression.* Let  $Z_i = (X_i, Y_i)$  with  $X_i \in [0, 1]^d$  and  $Y_i \in \mathbb{R}$  obeying the regression equation  $Y_i = f(X_i) + \xi_i$  with a regression function  $f$  and i.i.d. Gaussian errors  $\xi_i \sim \mathcal{N}(0, 1)$ .
2. *Inhomogeneous Poisson regression.* Here the observation is discrete  $Y_i \in \mathbb{N}$  and  $X_i \in [0, 1]^d$  is deterministic. The law of probability of  $Y_i$  can be written

$$g(k, f_i) = \mathbb{P}_f(Y_i = k) = \frac{[f(X_i)]^k}{k!} \exp\{-f(X_i)\}, \quad k \in \mathbb{N}.$$

This is the Poisson distribution with parameter  $f(X_i)$ . The goal is to estimate the function  $f$  at a given point  $y$ .

3.  *$\alpha$  Regression* Let  $0 < \alpha < 1/2$  and  $Z_i = (X_i, Y_i)$  satisfied the regression equation  $Y_i = f(X_i) + \epsilon_i$  where  $(\epsilon_i)_i$  are i.i.d. random variables of density  $g(x) = C(\alpha) e^{-|x|^\alpha}$ ,  $x \in \mathbb{R}$ .
4. *Multiplicative uniform regression.* Here  $Z_i = (X_i, Y_i)$  satisfied the regression equation  $Y_i = f(X_i) \times U_i$  where  $U_i \sim \mathcal{U}_{[0,1]}$ .

Along this thesis, the unknown function  $f$  is supposed to be smooth, in particular, it belongs to the Hölder ball of isotropic functions  $\mathfrak{F} = \mathbb{H}_d(\beta, L, M)$  (see Definition 11 below). Here  $\beta > 0$  is the smoothness of  $f$ ,  $M$  is the upper bound of  $f$  and its partial derivatives and  $L > 0$  is Lipschitz constant. We assume that we know *the minimax rate of convergence* (See Definition 12) given by the sequence

$$(3.1.2) \quad \varphi_{n,\gamma}(\beta) = n^{-\frac{\beta}{\gamma\beta+d}},$$

where the parameter  $\gamma$  is known and depend on the function  $g(\cdot)$ . For example, if  $\mathcal{Z}_n$  satisfied the *gaussian regression* then  $\gamma = 2$  and if  $\mathcal{Z}_n$  satisfied the *multiplicative uniform regression* then  $\gamma = 1$ .

**Minimax estimation.** The first part of this chapter is devoted to the minimax over  $\mathbb{H}_d(\beta, L, M)$  estimation. This means, in particular, that the parameters  $\beta, L$  and  $M$  are supposed to be known *a priori*. We propose the estimator being optimal in minimax sense (see Definition 12), i.e. the estimator attaining the normalising sequence (3.1.2). To construct the minimax estimator we use so-called *locally bayesian estimation construction* which consists in the following. Let

$$V_h(y) = \left\{ i = 1, \dots, n : X_i \in [0, 1]^d \cap \bigotimes_{j=1}^d [y_j - h/2, y_j + h/2] \right\},$$

be the neighborhood around  $y$ , where  $h \in (0, 1)$  be a given scalar. Fix  $b > 0$  (without loss of generality we will assume that  $b$  is integer) and let

$$D_b = \sum_{m=0}^b \binom{m+d-1}{d-1}.$$

Let  $K(z), z \in \mathbb{R}^d$  be the  $D_b$ -dimensional vector of polynomials of the following type (the sign  $\top$  below means the transposition):

$$K^\top(z) = \left( \prod_{j=1}^d z_j^{p_j}, (p_1, \dots, p_d) \in \mathbb{N}^d : 0 \leq p_1 + \dots + p_d \leq b \right).$$

For any  $t \in \mathbb{R}^{D_b}$  we define the local polynomial

$$(3.1.3) \quad f_t(x) = t^\top K \left( \frac{x-y}{h} \right) \mathbb{I}_{V_h(y)}(x), \quad x \in [0, 1]^d,$$

where  $\mathbb{I}$  denotes the indicator function. The local polynomial  $f_t$  can be viewed as an approximation of the regression function  $f$  inside of the neighborhood  $V_y(h)$ . Introduce the following subset of  $\mathbb{R}^{D_b}$

$$(3.1.4) \quad \Theta(M) = \{t \in \mathbb{R}^{D_b} : \|t\|_1 \leq M\},$$

where  $\|\cdot\|_1$  is  $\ell_1$ -norm on  $\mathbb{R}^{D_b}$ . For any  $i = 1, \dots, n$  remark that  $g(\cdot, f(X_i))$  is the density of the observation  $Y_i$ . Let  $\mathbb{E}_f = \mathbb{E}_f^n$  be the mathematical expectation with respect to the probability law  $\mathbb{P}_f = \mathbb{P}_f^n$  of the observation  $\mathcal{Z}_n$ . Consider the *pseudo likelihood ratio*

$$(3.1.5) \quad L_h(t, \mathcal{Z}_n) = \prod_{X_i \in V_h(y)} g(Y_i, f_t(X_i)), \quad t \in \Theta(M),$$

Set also

$$(3.1.6) \quad \pi_h(t) = \int_{\Theta(M)} \|t - u\|_1 [L_h(u, \mathcal{Z}_n)]^{1/m} du, \quad t \in \Theta(M).$$

where  $m$  is a positive constant chosen in Assumptions 1. Let  $\hat{\theta}(h)$  be the solution of the following minimization problem:

$$(3.1.7) \quad \hat{\theta}(h) = \arg \min_{t \in \Theta(M)} \pi_h(t).$$

The *locally bayesian estimator*  $\bar{f}^h(y)$  of  $f(y)$  is defined now as  $\bar{f}^h(y) = \hat{\theta}_{0, \dots, 0}(h)$ .

We note that similar locally parametric approach based on maximum likelihood estimators was recently proposed in [Polzehl and Spokoiny \[2006\]](#) and [Katzkovnik and Spokoiny \[2008\]](#) for *regular statistical models*.



As we see our construction contains an extra-parameter  $h$  to be chosen. To make this choice, we use quite standard arguments. First, we note that in view of  $f \in \mathbb{H}_d(\beta, L, M)$

$$\exists \theta = \theta(f, y, h) \in \Theta(M) : \sup_{x \in V_h(y)} |f(x) - f_\theta(x)| \leq Ldh^\beta.$$

Thus, if  $h$  is chosen sufficiently small our original observation (3.1.1) is well approximated inside of  $V_h(y)$  by the “parametric” observation  $\mathcal{Y}_i$  with the density  $g(\cdot - f_\theta(X_i))$  in which the *bayesian estimator*  $\hat{\theta}$  is rate-optimal (See Theorem 4).

Finally,  $h_n(\beta, L)$  is chosen as the solution of the following minimization problem

$$(3.1.8) \quad N_h^{-1} + Ldh^\beta \rightarrow \min_h$$

and we show that corresponding estimator  $\bar{f}^{h_n(\beta, L)}(y)$  is minimax for  $f(y)$  on  $\mathbb{H}_d(\beta, L, M)$  if  $\beta \leq b$ . Since the parameter  $b > 0$  can be chosen on arbitrary way, the proposed estimator is minimax for any given value of the parameter  $\beta > 0$ .

**Adaptive estimation.** The second part of this chapter is devoted to the adaptive minimax estimation over collection of isotropic functional classes. To our knowledge, the problem of optimal adaptive estimation in the general regression, with no-linear estimators, is not studied in the literature.

Well-known drawback of minimax approach is the dependence of the minimax estimator on the parameters describing functional class on which the maximal risk is determined. In particular, the locally bayesian estimator  $\bar{f}^h(\cdot)$  depends obviously on the parameter  $M$  via the solution of the minimization problem (3.1.7). Moreover,  $h_n(\beta, L)$  optimally chosen in view of (3.1.8) depends explicitly on  $\beta$  and  $L$ . To overcome this drawback the minimax adaptive approach was proposed (see Lepski [1990], Lepski [1991], Lepski, Mammen, and Spokoiny [1997]). The first question arising in the adaptation (reduced to the problem at hand) can be formulated as follows.

*Does there exist an estimator which would be minimax on  $\mathbb{H}_d(\beta, L, M)$  simultaneously for all values of  $\beta, L$  and  $M$  belonging to some given subset of  $\mathbb{R}_+^3$ ?*

In section 3.3, we determine that the answer on this question is **negative**, that is typical for the estimation of the function at a given point (Lepski and Spokoiny [1997]). This answer can be reformulated in the following manner: the family of rates of convergence  $\{\varphi_{n,\gamma}(\beta), \beta \in \mathbb{R}_+^*\}$  is **unattainable** for the problem under consideration.

Thus, we need to find another family of normalizations for maximal risk which would be attainable and, moreover, optimal in view of some criterion of optimality. Nowadays, the most developed criterion of optimality is due to Klutchnikoff [2005].

We can show that the family of normalizations, being optimal in several models (gaussian

$\gamma = 2$ , uniform  $\gamma = 1$  or  $\alpha$  regression  $\gamma = 1 + 2\alpha$ ) in view of this criterion, is

$$(3.1.9) \quad \phi_{n,\gamma}(\beta) = \left( \frac{\rho_{n,\gamma}(\beta)}{n} \right)^{\frac{\beta}{\gamma\beta+d}}, \quad \rho_{n,\gamma}(\beta) = \left[ 1 + \frac{\gamma(b-\beta)}{(\gamma b+d)(\gamma\beta+d)} \ln n \right]^{\frac{1}{\gamma}},$$

whenever  $\beta \in ]0, b]$ . The factor  $\rho_{n,\gamma}$  can be considered as *price to pay for adaptation* (Lepski [1990]).

The most important step in proving the optimality of the family (3.1.9) is to find an estimator, called *adaptive*, which attains the optimal family of normalizations. Obviously, we seek an estimator whose construction is *parameter-free*, i.e. independent of  $\beta$  and  $L$ . In order to explain our estimation procedure, let us make several remarks.

First, we note that the role of the constants  $\beta, L$  in the construction of the minimax estimator is quite different. Indeed, the constant  $M$  is used in order to determine the set  $\Theta(M)$  needed for the construction of the locally bayesian estimator, see (3.1.6) and (3.1.7). However, this set does not depend on the localisation parameter  $h > 0$ , in other words, the quantity  $M$  is not involved in the selection of optimal size of the local neighborhood given by (3.1.8). Contrary to that, the constants  $\beta, L$  are used for the derivation of the optimal size of the local neighborhood (3.1.8), but they are not involved in the construction of the collection of locally bayesian estimators  $\{\hat{f}^h, h > 0\}$ .

In order to select an “optimal” estimator from the family  $\{\hat{f}^h, h > 0\}$  we use general adaptation scheme due to Lepski [1990, 1992a]. Lepski’s procedure is typically applied to the selection from the collection of linear estimator (kernel estimators, locally polynomial estimator, etc.). In the present paper, we apply this method to a very complicated family of non-linear estimators, obtained by the use of bayesian approach on the random parameter set. It required, in particular, to establish the exponential inequality for the deviation of locally bayesian estimator from the parameter to be estimated (Proposition 3). It generalizes the inequality proved for the parametric model (Has’minskii and Ibragimov [1981]). This results seems to be new.

## 3.2 Minimax Estimation

In this section we present a upper bound concerning minimax estimation. First, we establish lower bound for minimax risk defined on  $\mathbb{H}_d(\beta, L, M)$  for any  $\beta, L$  and  $M$ . For any  $(p_1, \dots, p_d) \in \mathbb{N}^d$  we denote  $\vec{p} = (p_1, \dots, p_d)$  and  $|\vec{p}| = p_1 + \dots + p_d$ .

**Definition 11.** Fix  $\beta > 0$ ,  $L > 0$  and  $M > 0$  and let  $\lfloor \beta \rfloor$  be the largest integer strictly less than  $\beta$ . The isotropic Hölder class  $\mathbb{H}_d(\beta, L, M)$  is the set of functions  $f : [0, 1]^d \rightarrow \mathbb{R}$  having on  $[0, 1]^d$  all partial derivatives of order  $\lfloor \beta \rfloor$  and such that  $\forall x, y \in [0, 1]^d$

$$\sum_{m=0}^{\lfloor \beta \rfloor} \sum_{|\vec{p}|=m} \sup_{x \in [0,1]^d} \left| \frac{\partial^{|\vec{p}|} f(x)}{\partial x_1^{p_1} \cdots \partial x_d^{p_d}} \right| \leq M,$$

$$\left| \frac{\partial^{|\vec{p}|} f(x)}{\partial x_1^{p_1} \cdots \partial x_d^{p_d}} - \frac{\partial^{|\vec{p}|} f(y)}{\partial y_1^{p_1} \cdots \partial y_d^{p_d}} \right| \leq L [\|x - y\|_1]^{\beta - \lfloor \beta \rfloor}, \quad \forall |\vec{p}| = \lfloor \beta \rfloor.$$

**Maximal and Minimax Risk on  $\mathbb{H}_d(\beta, L, M)$ .** To measure the performance of estimation procedures on  $\mathbb{H}_d(\beta, L, M)$  we will use minimax approach.

First, we define the maximal risk on  $\mathbb{H}_d(\beta, L, M)$  corresponding to the estimation of the function  $f$  at a given point  $y \in [0, 1]^d$ .

Let  $\tilde{f}$  be an arbitrary estimator built from the observation  $\mathcal{Z}_n$ . Let  $\forall q > 0$

$$R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M)] = \sup_{f \in \mathbb{H}_d(\beta, L, M)} \mathbb{E}_f |\tilde{f}(y) - f(y)|^q.$$

The quantity  $R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M)]$  is called *maximal risk* of the estimator  $\tilde{f}$  on  $\mathbb{H}_d(\beta, L, M)$  and the *minimax risk* on  $\mathbb{H}_d(\beta, L, M)$  is defined as

$$R_{n,q}[\mathbb{H}_d(\beta, L, M)] = \inf_{\tilde{f}} R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M)],$$

where inf is taken over the set of all estimators.

**Definition 12.** The normalising sequence  $\psi_n$  is called *minimax rate of convergence* and the estimator  $\hat{f}$  is called *minimax (asymptotically minimax)* if

$$\liminf_{n \rightarrow \infty} \psi_n^{-q} R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M, A)] > 0;$$

$$\limsup_{n \rightarrow \infty} \psi_n^{-q} R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M, A)] < \infty.$$

We remember that the first assumption of this chapter is that  $\varphi_{n,\gamma}$  defined in (3.1.2) is the *minimax rate of convergence*. To obtain the upper bound of minimax risk, we need to the following assumptions.

Let us introduce the following notations. For any  $\theta \in \Theta(M)$ , let  $U_n = N_h [\Theta(M) - \theta]$ . Put also  $N_h = N_h(\gamma) = (nh^d)^{1/\gamma}$  where  $\gamma \geq 1$  and  $\forall u \in U_n$

$$(3.2.1) \quad Z_{n,\theta}(u) = \frac{L_h(\theta + u N_h^{-1}, \mathcal{Z}_n)}{L_h(\theta, \mathcal{Z}_n)}.$$

Let  $\mathcal{H}_n, n > 1$  be the following subinterval of  $(0, 1)$ .

$$(3.2.2) \quad \mathcal{H}_n = [h_{\min}, h_{\max}], \quad h_{\min} = (\ln n)^{\frac{1}{\gamma d + d^2}} n^{-1/d}, \quad h_{\max} = (\ln n)^{-\frac{1}{\gamma b + d}}.$$

Later on we will consider only the values of  $h$  belonging to  $\mathcal{H}_n$ . Let  $f_\theta(x)$ , with  $\theta \in \Theta(M)$ , be the local polynomial approximation of  $f$  inside  $V_h(y)$  and let  $b_h$  be the corresponding approximation error, i.e.

$$(3.2.3) \quad b_h = \sup_{x \in V_h(y)} |f_\theta(x) - f(x)|.$$

Put finally  $\mathcal{N}(h) = (b_h \times N_h)^\gamma$  and define

$$(3.2.4) \quad \mathcal{E}_h = \exp \{ \mathcal{N}(h) \}.$$

Let the subset  $\Gamma_\delta \subseteq \Theta(M)$  such that  $\int_{\Gamma_\delta} du = \delta^{D_b}$ ,  $\delta > 0$ .

**Assumptions 1.** *We suppose that it exists several constants  $\tau, m > 0$  and  $c_1, c_2, c_3, C_1, C_2, s_1, s_2 > 0$  such that for any  $\forall n > 1$ ,  $f \in \mathbb{H}_d(\beta, L, M)$ ,  $\theta \in \Theta(M)$  and  $h \in \mathcal{H}_n$ , we have*

1.  $\int_{\Gamma_\delta} \mathbb{E}_f \left[ 1 - Z_{n,\theta}^{1/m}(u) \right]_+ du \leq C_1 \mathcal{E}_h^{c_1} \delta^{D_b + \tau}, \quad \forall \delta < s_1,$
2.  $\mathbb{E}_f Z_{n,\theta}^{1/m}(u) \leq C_2 \mathcal{E}_h^{c_2} \exp \{ -c_3 \|u\|_1^\gamma \}, \quad \forall u \in U_n : \|u\|_1 \geq s_2 \mathcal{N}(h).$

where  $a_+ = \max(a, 0)$ ,  $a \in \mathbb{R}$ . We also assume the acknowledgement of constants  $c_3, s_2, \tau, m$  and  $\gamma$ .

The next theorem shows how to construct the minimax estimator basing on locally bayesian approach. Put  $\bar{h} = (L^\gamma n)^{-\frac{1}{\gamma\beta+d}}$  and let  $\bar{f}^h(y) = \hat{\theta}_{0,\dots,0}(\bar{h})$  is given by (3.1.4), (3.1.6) and (3.1.7) with  $h = \bar{h}$ .

**Theorem 4.** *Let  $\beta > 0$ ,  $L > 0$  and  $M > 0$  be fixed. Assume that Assumptions 1 are verified. Then it exists the constant  $C_*$  such that for any  $n \in \mathbb{N}^*$  satisfying  $N_{\bar{h}} \geq 1$ ,*

$$\varphi_{n,\gamma}^{-q}(\beta) R_{n,q} \left[ \bar{f}^h(y), \mathbb{H}_d(\beta, L, M) \right] \leq C_*, \quad \forall q \geq 1.$$

### 3.3 Adaptive Rule

This section is devoted to the adaptive estimation over the collection of the classes  $\left\{ \mathbb{H}_d(\beta, L, M) \right\}_{\beta, L, M}$ . We will not impose any restriction on possible values of  $L, M$ , but we will assume that  $\beta \in (0, b]$ , where  $b$ , as previously, is an arbitrary *a priori* chosen integer.

Let  $\Phi$  be the following family of normalizations:

$$\phi_{n,\gamma}(\beta) = \left( \frac{\rho_{n,\gamma}(\beta)}{n} \right)^{\frac{\beta}{\gamma\beta+d}}, \quad \rho_{n,\gamma}(\beta) = \left[ 1 + \frac{\gamma(b-\beta)}{(\gamma b + d)(\gamma\beta + d)} \ln n \right]^{\frac{1}{\gamma}}, \quad \beta \in (0, b].$$

We remark that  $\phi_{n,\gamma}(b) = \varphi_{n,\gamma}(b)$  and  $\rho_{n,\gamma}(\beta) \sim (\ln n)^{\frac{1}{\gamma}}$  for any  $\beta \neq b$ .

**Construction of  $\Phi$ -Adaptive Estimator.** As it was already mentioned in Introduction the construction of our estimation procedure consists of two steps. First, we define the family of locally bayesian estimators. Next, based on Lepski's method, we propose data-driven selection from this family.

First step: Collection of locally bayesian estimators. Put

$$h_k = 2^{-k} h_{\max}, \quad k = 0, \dots, \mathbf{k}_n,$$

where  $\mathbf{k}_n$  is largest integer such that  $h_{\mathbf{k}_n} \in \mathcal{H}_n$  defined in (3.2.2). Set

$$\hat{\mathcal{F}} = \left\{ \hat{f}^{(k)}(y) = \hat{\theta}_{0,\dots,0}(h_k), \quad k = 0, \dots, \mathbf{k}_n \right\},$$

where  $\hat{\theta}_{0,\dots,0}(h_k)$  is given by (3.1.4), (3.1.6) and (3.1.7) with  $h = h_k$ .

Second step: Data-driven selection from the collection  $\hat{\mathcal{F}}$ . We put  $\hat{f}^*(y) = \hat{f}^{(\hat{k})}(y)$ , where  $\hat{f}^{(\hat{k})}(y)$  is selected from  $\hat{\mathcal{F}}$  in accordance with the rule:

$$(3.3.1) \quad \hat{k} = \inf \left\{ k = \overline{0, \mathbf{k}_n} : |\hat{f}^{(k)}(y) - \hat{f}^{(l)}(y)| \leq C S_n(l), \quad l = \overline{k+1, \mathbf{k}_n} \right\}.$$

Here we have used the following notation.

$$(3.3.2) \quad C = \left( s_2^\gamma \vee \frac{2^{2\gamma+1}(2\gamma + 2dq)}{c_3\gamma(1 \wedge \tau D_b^{-1})} \right)^{\frac{1}{\gamma}}, \quad S_n(l) = \left[ \frac{1 + l \ln 2}{n(h_l)^d} \right]^{\frac{1}{\gamma}}, \quad l = 0, 1, \dots, k_n.$$

**Theorem 5.** Let  $b > 0$  be fixed. Assume that Assumptions 1 are verified, then for any  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$  and  $q \geq 1$

$$\limsup_{n \rightarrow \infty} \phi_{n,\gamma}^{-q}(\beta) R_{n,q} \left[ \hat{f}^*(y), \mathbb{H}_d(\beta, L, M) \right] < \infty.$$

**Remark 7.** The assertion of the theorem means that the proposed estimator  $\hat{f}^*(y)$  is  $\Phi$ -adaptive. It implies in particular that the family of normalizations  $\Phi$  is admissible. For example (Section 4.6), we can state the optimality of  $\Phi$  in view of Klutchnikoff criterion.

## 3.4 Applications

In this section, we show how the upper bounds of Theorems 4 and 5 can be used for the derivation of minimax and adaptive minimax results for different additive models of type 3.1.1. In particular, in following sections, we consider four particular problems:

1.  $g(\cdot - f_i)$  is the normal distribution  $\mathcal{N}(0, 1)$ ,

2.  $g(k, f_i) = \mathbb{P}_f(Y_i = k) = \frac{[f(X_i)]^k}{k!} \exp\{-f(X_i)\}$ ,  $k \in \mathbb{N}$ ,
3.  $g(x, f_i) = C(\alpha) \exp\{-|x - f_i|^\alpha\}$ ,  $x \in \mathbb{R}$ ,
4.  $g(\cdot, f_i)$  is the uniform distribution on  $[0, f_i]$ .

Our goal here is to show how the assumption on  $g(\cdot)$  leads to bayesian estimators with optimal statistical properties. Note that in each particular case for adaptation, the bayesian estimators are obtained by the same selection rule presented in Section 3.3. Here, we always assume that the unknown  $f \in \mathbb{H}_d(\beta, L, M)$ .

### 3.4.1 Gaussian Regression

Consider the model 3.1.1 with  $g(x, f_i) = (\sqrt{2\pi})^{-1} \exp\{-(x - f_i)^2/2\}$ , we obtain the following model:

$$(3.4.1) \quad Y_i = f(X_i) + \xi_i, \quad i = 1, \dots, n.$$

According to Tsybakov [2008], it is well-known that the minimax rate is  $n^{-\frac{\beta}{2\beta+d}}$ , thus here we take  $\gamma = 2$ . By definition (3.2.1) the process  $Z_{n,\theta}$  can be formulate like

$$Z_{n,\theta}(u) = \prod_{X_i \in V_y(h)} \exp\left\{-2^{-1}(Y_i - f_{\theta+uN_h^{-1}}(X_i))^2 + 2^{-1}(Y_i - f_\theta(X_i))^2\right\}$$

Put  $\lambda_n(h)$  the smallest eigenvalue of the matrix

$$(3.4.2) \quad \mathcal{M}_{nh}(y) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{X_i - y}{h}\right) K^\top\left(\frac{X_i - y}{h}\right) \mathbb{I}_{V_h(y)}(X_i),$$

which is completely determined by the design points and by the number of observations. We will prove that it exists  $\lambda > 0$ , such that  $\lambda_n(h) \geq \lambda$  for any  $n \geq 1$  and any  $h \in [h_{\min}, h_{\max}]$  (see Lemma 9 in Chapter 4). Let us give the following lemma which verified Assumptions 1 on  $Z_{n,\theta}$ .

**Lemma 1.** *For any  $f \in \mathbb{H}_d(\beta, L, M)$ ,  $n \geq 1$ ,  $h \in \mathcal{H}_n$ ,  $\theta \in \Theta(M)$*

1.  $\int_{[0,\delta]^{D_b}} \mathbb{E}_f \left[1 - Z_{n,\theta}^{1/2}(u)\right]_+ du \leq 2\mathcal{E}_h \delta^{D_b+2}, \quad \forall \delta < 1,$
2.  $\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq \exp\left\{-\frac{D_b \lambda_n(h)}{16} \|u\|_1^2\right\}, \quad \forall u \in U_n : \|u\|_1 \geq \frac{8}{D_b \lambda_n(h)} \mathcal{N}(h).$

The proof is given in Appendix. We use the adaptive bayesian estimator noted  $\hat{f}_1^*(y)$  and developed in Section 3.3 with model 3.4.1, by identification we take constants  $\gamma = 2$ ,  $\tau = 2$ ,  $m = 2$ ,  $s_2 = 8 D_b^{-1} \lambda_n^{-1}(h)$ ,  $c_3 = 2^{-4} D_b \lambda_n(h)$ .

**Theorem 6.** *Let  $b > 0$  be fixed, then for any  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$  and  $q \geq 1$*

$$\limsup_{n \rightarrow \infty} \phi_{n,2}^{-q}(\beta) R_{n,q} \left[ \hat{f}_1^*(y), \mathbb{H}_d(\beta, L, M) \right] < \infty.$$

**Remark 8.** *The proof of the theorem is omitted; it is a straightforward consequence of Theorem 5 and Lemma 1. We can deduced the minimax result to take  $b = \beta$  and remark that  $\phi_{n,2}(b) = \varphi_{n,2}(b)$ . With Klutchnikoff criterion, we can prove that  $\{\phi_{n,2}(\beta)\}_{\beta \in [0,b]}$  is adaptive optimal for the Gaussian regression and the pointwise estimation (See [Klutchnikoff \[2005\]](#)).*

### 3.4.2 Inhomogeneous Poisson Regression

Consider the model 3.1.1 with  $g(k, f_i) = \mathbb{P}_f(Y_i = k) = \frac{[f(X_i)]^k}{k!} \exp\{-f(X_i)\}$ ,  $k \in \mathbb{N}$ . Assume that  $f \in \mathbb{H}(\beta, L, M, A)$ ,  $M \geq A > 0$ , where

$$\mathbb{H}(\beta, L, M, A) = \left\{ f \in \mathbb{H}(\beta, L, M) : \inf_{x \in [0,1]^d} f(x) \geq A \right\}.$$

By definition (3.2.1) the process  $Z_{n,\theta}$  can be formulate like

$$Z_{n,\theta}(u) = \prod_{X_i \in V_{y(h)}} \left( 1 + \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^{Y_i} \exp \left\{ -f_{uN_h^{-1}}(X_i) \right\}, \quad N_h = (nh^d)^{1/2}.$$

Put  $\lambda_n(h)$  given in the previous section. Let us give the following lemma which verified Assumptions 1 on  $Z_{n,\theta}$ . Let us introduce the following notation

$$(3.4.3) \quad \Theta(A, M) = \{t \in \mathbb{R}^{D_b} : 2t_{0,\dots,0} - \|t\|_1 \geq A, \quad \|t\|_1 \leq M\}.$$

**Lemma 2.** *For any  $f \in \mathbb{H}_d(\beta, L, M)$ ,  $n \geq n_0$ ,  $h \in \mathcal{H}_n$ ,  $\theta \in \Theta(A, M)$*

1.  $\int_{[0,\delta]^{D_b}} \mathbb{E}_f \left[ 1 - Z_{n,\theta}^{1/2}(u) \right]_+ du \leq \sqrt{\frac{2+M+2A}{2A}} \mathcal{E}_h^{1/2} \delta^{D_b+1/2}, \quad \forall \delta < 1,$
2.  $\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq e^1 \exp \left\{ -\frac{\lambda_n(h)A}{8M^2} \|u\|_1^2 \right\}, \quad \forall u \in U_n : \|u\|_1 \geq \frac{4M^2}{\lambda_n(h)} \mathcal{N}(h).$

Here,  $U_n = N_h(\Theta(A, M) - \theta)$ .

The proof is given in Appendix. Remark that this result is given for  $n \geq n_0$ , this do not change the final result which is asymptotic. We use the adaptive bayesian estimator noted  $\hat{f}_2^*(y)$  and developed in Section 3.3, by identification we take constants  $\gamma = 2$ ,  $\tau = 0.5$ ,  $m = 2$ ,  $s_2 = 4M^2\lambda_n^{-1}(h)$ ,  $c_3 = 2^{-3}M^{-2}\lambda_n(h)A$ .

**Theorem 7.** *Let  $b > 0$  be fixed, then for any  $\beta \in (0, b]$ ,  $L > 0$ ,  $M \geq A > 0$  and  $q \geq 1$*

$$\limsup_{n \rightarrow \infty} \phi_{n,2}^{-q}(\beta) R_{n,q} \left[ \hat{f}_2^*(y), \mathbb{H}_d(\beta, L, M, A) \right] < \infty.$$

**Remark 9.** *The proof of the theorem is omitted; it is a straightforward consequence of Theorem 5 and Lemma 2. We can deduced the minimax result to take  $b = \beta$  and remark that  $\phi_{n,2}(b) = \varphi_{n,2}(b)$ . A good exercice is to prove that  $\{\phi_{n,2}(\beta)\}_{\beta \in [0,b]}$  is adaptive optimal for the Poisson regression and the pointwise estimation (See Klutchnikoff [2005]). In this case, the constant  $c_3$  depends explicitly of  $A, M$ , unknown in practice. We can construct no-optimal estimators of  $A, M$  to choose the constant  $C$  in the Lepski's procedure. Few examples of estimators are given for the multiplicative uniform regression (Chapter 4).*

### 3.4.3 $\alpha$ Regression

We call  $\alpha$  regression the following model

$$(3.4.4) \quad Y_i = f(X_i) + \epsilon_i, \quad i = 1, \dots, n,$$

where  $\epsilon_i$  has for density  $g_\alpha(x) = C(\alpha) \exp\{-|x|^\alpha\}$  with  $0 < \alpha < 1/2$ .  $C(\alpha)$  is a constant chosen so that  $g_\alpha$  density is well. Here the design  $X_i$  is deterministic uniform on  $[0, 1]^d$  and  $f \in \mathbb{H}_d(\beta, L, M)$ .

In this model, we want to achieve the rate of convergence  $n^{-\frac{\beta}{(1+2\alpha)\beta+d}}$  (given by Has'minskii and Ibragimov [1981] in the parametric estimation). Thus, we take  $\gamma = 1 + 2\alpha$ . By definition (3.2.1) the process  $Z_{n,\theta}$  can be formulate like

$$Z_{n,\theta}(u) = \prod_{X_i \in V_y(h)} \exp \left\{ -|Y_i - f_{\theta+uN_h^{-1}}(X_i)|^\alpha + |Y_i - f_\theta(X_i)|^\alpha \right\}$$

Put  $\lambda_n(h)$  given in Section 3.4.1. Let us give the following lemma which verified Assumptions 1 on  $Z_{n,\theta}$ .

**Lemma 3.** *For any  $f \in \mathbb{H}_d(\beta, L, M)$ , it exists  $c_1, \mathcal{C}, n_0$  such that for  $n \geq n_0$ ,  $h \in \mathcal{H}_n$ ,  $\theta \in \Theta(M)$ , we have*

$$1. \int_{[0,\delta]^{D_b}} \mathbb{E}_f \left[ 1 - Z_{n,\theta}^{1/2}(u) \right]_+ du \leq \mathcal{C} \mathcal{E}_h^{c_1} \delta^{D_b+\alpha+1/2}, \quad \forall \delta < 1,$$

$$2. \mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq \mathcal{E}_h^{\frac{1+C(\alpha)}{2(1+\alpha)}} \exp \left\{ -c_3 \|u\|_1^{1+2\alpha} \right\}, \quad \forall u \in U_n,$$

$$\text{where } c_3 = \lambda_n(h) 2^{\alpha-4} C(\alpha) (1 + 2\alpha)^{-1} \exp \left\{ -\frac{1}{2} (2M)^\alpha \right\}.$$

The proof is given in Appendix. Remark that we do not know values of  $c_1, \mathcal{C}$ , it is not necessary to construct the adaptive estimator. We use the adaptive bayesian estimator noted  $\hat{f}_3^*(y)$  and developed in Section 3.3 with model 3.4.4.



**Theorem 8.** *Let  $b > 0$  be fixed, then for any  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$  and  $q \geq 1$*

$$\limsup_{n \rightarrow \infty} \phi_{n,1+2\alpha}^{-q}(\beta) R_{n,q} \left[ \hat{f}_3^*(y), \mathbb{H}_d(\beta, L, M) \right] < \infty.$$

**Remark 10.** *The proof of the theorem is omitted, it is a straightforward consequence of Theorem 5 and Lemma 3. The optimality of  $\phi_{n,1+2\alpha}(\beta)$  is not proved, it is an open problem. Here Lemma 3 is given with  $n \geq n_0$ , it will not be a problem because Theorem 8 is an asymptotic result. We can construct the minimax estimator with Theorem 4, because in the case minimax, the term  $\mathcal{E}_h$  (in Assumptions 1) is a constant for  $h = \bar{h}$ .*

### 3.4.4 Multiplicative Uniform Regression

We consider the following model:

$$(3.4.5) \quad Y_i = f(X_i) \times U_i, \quad i = 1, \dots, n,$$

where  $(U_i)_i$  is a sequence independent identically distributed of uniform law on  $[0, 1]$ . This chapter 4 is devoted to the study of this model. So, we can find the minimax rate  $n^{-\frac{\beta}{\beta+d}}$ , thus, here we take  $\gamma = 1$ . The proof of Assumptions 1 is a consequence of Lemma 8 in Chapter 4 and we give the following lemma

**Lemma 4.** *For any  $f \in \mathbb{H}_d(\beta, L, M, A)$ ,  $n \geq 1$ ,  $h \in \mathcal{H}_n$ ,  $\theta \in \Theta(A, M)$*

1.  $\int_{[0,\delta]^{D_b}} \mathbb{E}_f \left[ 1 - Z_{n,\theta}^{1/2}(u) \right]_+ du \leq \frac{3}{2A} \delta^{D_b+1}, \quad \forall \delta < 1,$
2.  $\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq \mathcal{E}_h^{1/A} \exp \left\{ -\frac{\lambda_n(h)}{18MD_b} \|u\|_1 \right\}, \quad \forall u \in N_h(\Theta(A, M) - \theta).$

where

$$Z_{n,\theta}(u) = \prod_{X_i \in V_y(h)} \frac{f_\theta(X_i)}{f_{\theta+uN_h^{-1}}(X_i)} \mathbb{I}_{\{Y_i \leq f_{\theta+uN_h^{-1}}(X_i)\}}.$$

This case requires us to several restrictions on the function  $f$ . Indeed, we assume that  $f \geq A > 0$  ( $f \in \mathbb{H}_d(\beta, L, M, A)$ ) and the polynomial approximation is constructed with the constraint  $f_\theta \geq f$ . We find a lot of details in Chapter 4, in particular this lemma is a consequence of Lemma 8.

We use the adaptive bayesian estimator noted  $\hat{f}_4^*(y)$  and developed in Section 3.3 with model 3.4.5. By identification, we take constants  $\gamma = 1$ ,  $\tau = 1$ ,  $m = 2$ ,  $s_2 = 0$ ,  $c_3 = \frac{\lambda_n(h)}{18MD_b}$ .

**Theorem 9.** *Let  $b > 0$  be fixed, then for any  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$  and  $q \geq 1$*

$$\limsup_{n \rightarrow \infty} \phi_{n,1}^{-q}(\beta) R_{n,q} \left[ \hat{f}_4^*(y), \mathbb{H}_d(\beta, L, M, A) \right] < \infty.$$

**Remark 11.** *The proof of the theorem is omitted, it is a straightforward consequence of Theorem 5 and Lemma 4. With Klutchnikoff criterion, we can prove that  $\{\phi_{n,1}(\beta)\}_{\beta \in [0,b]}$  is adaptive optimal for the multiplicative uniform regression and the pointwise estimation (See Theorem 13 in Chapter 4).*

## 3.5 Proofs of Main Results

In this section we give the proofs of Theorems 4 and 5.

### 3.5.1 Auxiliary Results: Large Deviations

For the sequel we fixe, for any  $h > 0$  satisfying (3.2.2) and  $f \in \mathbb{H}_d(\beta, L, M)$ ,  $\theta = \theta(f, y, h) = \{\theta_{\vec{p}} : \vec{p} \in \mathbb{N}^d : 0 \leq |\vec{p}| \leq b\}$  such that  $b_h \leq Ldh^\beta$  and  $\theta \in \Theta(M)$ . For example, we can take  $\theta$  the coefficients of expansion Taylor series.

The next proposition is the milestone for all results proved in the paper. Introduce the following notations  $\omega_2 = c_1 2^{-1} \vee c_2$ ,

$$\omega_1 = 16 C_2 \int_0^{+\infty} (z+1) e^{-c_3 z^\gamma} dz + 4\sqrt{C_1}, \quad \omega_3 = c_3 2^{-2\gamma-1} \left(1 \wedge \frac{\tau}{D_b}\right),$$

where  $C_1, C_2, c_1, c_2, c_3$  are defined in section 3.2. Put

$$(3.5.1) \quad M(f) = \sum_{m=0}^b \sum_{p_1+\dots+p_d=m} \left| \frac{\partial^m f(y)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} \right|.$$

The following agreement will be used in the sequel: if the function  $f$  and  $m \geq 1$  be such that  $\partial^m f$  does not exist we will put formally  $\partial^m f = 0$  in the definition of  $M(f)$ .

**Proposition 3.** *Let  $\varepsilon_0 = 4 s_2 \vee \left(c_3^{-1} 2^{2\gamma+1} (M \wedge \ln s_1)_+\right)^{1/\gamma}$ . For any  $n > 1$ ,  $h \in \mathcal{H}_n$  and any  $f$  such that  $M(f) < \infty$ , then  $\forall \varepsilon \geq \varepsilon_0 \mathcal{N}(h)$*

$$\mathbb{P}_f \left( N_h |\hat{f}^h(y) - f(y)| \geq \varepsilon \right) \leq \omega_1 \mathcal{E}_h^{\omega_2} \exp \{-\omega_3 \varepsilon^\gamma\},$$

where  $\hat{f}^h(y)$  is given by (3.1.4), (3.1.6) and (3.1.7) and  $s_1, s_2$  are defined in Assumptions 1 in Section 3.2.

The proof of this proposition is given in section 3.5.4.

### 3.5.2 Proof of Theorem 4

By integration of the inequality of Proposition 3, we obtain for any  $n > 1$ ,  $h \in \mathcal{H}_n$  and any  $f$  such that  $M(f) < \infty$

$$(3.5.2) \quad \mathbb{E}_f |\hat{f}^h(y) - f(y)|^q \leq C_q \mathcal{E}_h^{[q \vee (\omega_2 - \omega_3 \varepsilon_0/2)]} N_h^{-q},$$

where  $C_q = 2\varepsilon_0^q + \frac{4q\omega_1}{\omega_3} \sup_{\eta \geq 0} [\eta^{q-1} \exp\{-\omega_3 \eta^\gamma/2\}]$ . By definition of  $\bar{h} = (Ln)^{-\frac{1}{\beta+d}} \in \mathcal{H}_n$ , we have

$$Ldn h^{\gamma\beta+d} = d, \quad N_{\bar{h}}^{-q} = L^{\frac{\gamma q d}{\gamma\beta+d}} \varphi_{n,\gamma}^q(\beta).$$

Using (3.5.2) and last equalities, we obtain for any  $f \in \mathbb{H}_d(\beta, L, M)$

$$\mathbb{E}_f |\bar{f}^{\bar{h}}(y) - f(y)|^q \leq C_q e^{d[q \vee (\omega_2 - \omega_3 \varepsilon_0/2)]} L^{\frac{q d \gamma}{\gamma\beta+d}} \varphi_{n,\gamma}^q(\beta).$$

Theorem 4 is proved. ■

### 3.5.3 Proof of Theorem 5

We start the proof with formulating some auxiliary results whose proofs are given in Appendix 3.6. Define  $J_1 = \omega_2$  and

$$h^* = \left[ \frac{c \rho_{n,\gamma}^\gamma(\beta)}{L^\gamma d^\gamma n} \right]^{\frac{1}{\gamma\beta+d}}, \quad c < (2J_1)^{-\gamma}.$$

and let the integer  $\kappa$  be defined as follows.

$$(3.5.3) \quad 2^{-\kappa} h_{\max} \leq h^* < 2^{-\kappa+1} h_{\max}.$$

**Lemma 5.** *For any  $f \in \mathbb{H}_d(\beta, L, M)$  and any  $k \geq \kappa + 1$*

$$\mathbb{P}(\hat{k} = k) \leq J_2 \exp \left\{ J_1 n (h^*)^{\gamma\beta+d} \right\} 2^{-(k-1) \frac{2\gamma+2dq}{\gamma}},$$

where  $J_2 = \omega_1 \left( 1 - 2^{-\frac{2\gamma+2dq}{\gamma}} \right)^{-1}$ .

**Proof of Theorem 5.** Note that  $h^* \in \mathcal{H}_n$  and  $h_l \in \mathcal{H}_n$ ,  $l = 0, \dots, \mathbf{k}_n$ . The definition of  $h^*$  and  $\kappa$  in (3.5.3) implies that by integration of the inequality obtained in Proposition 3 with  $h = h_k$ , we obtain:

$$(3.5.4) \quad \mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^q \leq \bar{C}_q \frac{(1 + k \ln 2)^{\frac{q}{\gamma}}}{N_{h_k}^q}, \quad \forall k \geq \kappa,$$

where  $\bar{C}_q = 2 \left(1 \vee c^{1/\gamma}\right)^q \left(\frac{\omega_2}{\omega_3} \varepsilon_0 \vee 1\right)^q \left(1 + q \omega_1 \int_0^{+\infty} (\eta + 1)^{q-1} e^{-\omega_3 \eta^\gamma} d\eta\right)$ . We also have

$$\begin{aligned}
 & \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \\
 & \leq \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_{\hat{k} \leq \kappa} + \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_{\hat{k} > \kappa} \\
 (3.5.5) \quad & := R_1(f) + R_2(f).
 \end{aligned}$$

First, we control  $R_1$ . Obviously

$$|\hat{f}^{(\hat{k})}(y) - f(y)| \leq |\hat{f}^{(\hat{k})}(y) - \hat{f}^{(\kappa)}(y)| + |\hat{f}^{(\kappa)}(y) - f(y)|.$$

The definition of  $\hat{k}$  yields

$$|\hat{f}^{(\hat{k})}(y) - \hat{f}^{(\kappa)}(y)| \mathbb{I}_{\hat{k} \leq \kappa} \leq C s_n(\kappa), \quad s_n(k) = (1 + k \ln 2)^{\frac{q}{\gamma}} N_{h_k}^{-q},$$

where  $C = \left(\omega_3^{-1} \frac{4\gamma + 4dq}{\gamma}\right)^{1/\gamma}$ . In view of (3.5.4) we also get

$$\mathbb{E}_f |\hat{f}^{(\kappa)}(y) - f(y)|^q \leq \bar{C}_q s_n(\kappa).$$

Noting that the right hand side of the obtain inequality is independent of  $f$  and taking into account the definition of  $\kappa$  and  $h^*$  we obtain

$$(3.5.6) \quad \limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, M)} \phi_n^{-q}(\beta) R_1(f) < \infty.$$

Now let us bounded from above  $R_2$ . Applying Cauchy-Schwarz inequality we have in view of Lemma 5

$$\begin{aligned}
 R_2(f) &= \sum_{k > \kappa}^{\mathbf{k}_n} \mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^q I_{[\hat{k}=k, G]} \\
 &\leq \sum_{k > \kappa} (\mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^{2q})^{1/2} \sqrt{\mathbb{P}_f \{\hat{k} = k\}} \\
 (3.5.7) \quad &= \Delta(h^*) \sum_{k > \kappa} (\mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^{2q})^{1/2} 2^{-(k-1) \frac{\gamma+dq}{\gamma}},
 \end{aligned}$$

where we have put  $\Delta(h^*) = J_2 \exp \{J_1 n(h^*)^{\gamma\beta+d}\}$ . We obtain from (3.5.4) and (3.5.7)

$$(3.5.8) \quad R_2(f) \leq J_3 N_{h_{\max}}^{-q} \exp \{J_1 n(h^*)^{\gamma\beta+d}\},$$

where

$$J_3 = J_2 2^{\frac{\gamma+dq}{\gamma}} \bar{C}_{2q}^{1/2} \sum_{s \geq 0} (1 + s \ln 2)^{\frac{q}{\gamma}} 2^{-s}.$$

It remains to note that the definition of  $h^*$  implies that

$$\limsup_{n \rightarrow \infty} \phi_n^{-q}(\beta) N_{h_{\max}}^{-q} \exp \{J_1 n (h^*)^{\gamma\beta+d}\} < \infty$$

and that the right hand side of (3.5.8) is independent of  $f$ . Thus, we have

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, M)} \phi_n^{-q}(\beta) R_2(f) < \infty.$$

that yields together with (3.5.5) and (3.5.6)

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, M)} \phi_n^{-q}(\beta) \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q < \infty.$$

■

### 3.5.4 Proof of Proposition 3

**Auxiliary Lemmas.** Set for any  $a > 0$

$$Q_a = \int_{U_n \cap \{\|u\|_1 > a\}} \|u\|_1 Z_{n,\theta}^{1/m}(u) du.$$

where  $U_n := N_h(\Theta - \theta)$ ,  $\Theta$  defined in (3.1.4) and  $Z_{n,\theta}$  defined in (3.2.1). Remember that  $\gamma > 0$  is the constant which defined the *minimax rate of convergence*  $\varphi_{n,\gamma}$  and depends on the density  $g(\cdot)$ .

**Lemma 6.** *For any  $f$  such that  $M(f) < \infty$  and  $h \in \mathcal{H}_n$ , if the assumption (1) is verified, then*

$$\mathbb{P}_f \left\{ \int_{U_n} Z_{n,\theta}^{1/m}(u) du < \frac{1}{2} \delta^{D_b} \right\} < 2C_1 \mathcal{E}_h^{c_1} \delta^\tau, \quad \forall \delta \in ]0, 1 \wedge M].$$

**Lemma 7.** *For all  $h \in \mathcal{H}_n$  and any  $f$  such that  $M(f) < \infty$ , if the assumption 1.2 is verified, then for any  $a > s_2 \mathcal{N}(h)$*

$$\mathbb{E}_f \{Q_a\} \leq 2C_2 a \mathcal{E}_h^{c_2} I_1 e^{-c_3 a^\gamma},$$

where  $s_2$  is defined in Section 3.2 and

$$I_1 = \int_0^{+\infty} (z+1) e^{-c_3 z^\gamma} dz$$

This lemmas are proved in section 3.6.

**Proof of Proposition 3.** The definition of  $\hat{\theta}(h)$  and  $\theta = \theta(f, y, h)$  implies  $\forall \varepsilon > 0$

$$(3.5.9) \quad \begin{aligned} \mathbb{P}_f \left( N_h | \hat{f}^h(y) - f(y) | \geq \varepsilon \right) &\leq \mathbb{P}_f \left( N_h | \hat{\theta}_0(h) - \theta_0 | \geq \varepsilon \right) \\ &\leq \mathbb{P}_f \left( N_h \| \hat{\theta}(h) - \theta \|_1 \geq \varepsilon \right). \end{aligned}$$

Let us do some remarks. First, by definition we have  $\theta \in \Theta(M)$ . Remind also that  $\hat{\theta}(h)$  minimizes  $\pi_h$  defined in (3.1.6) and, therefore, the following inclusion holds since  $\hat{\theta}(h) \in \Theta$ .

$$(3.5.10) \quad \left\{ N_h \| \hat{\theta}(h) - \theta \|_1 \geq \varepsilon \right\} \subseteq \left\{ \inf_{N_h \| t - \theta \|_1 \geq \varepsilon} \pi_h(t) \leq \pi_h(\theta) \right\}.$$

Moreover,

$$\begin{aligned} \pi_h(t) &= N_h^{-1} \int_{\Theta} \| N_h(t - u) \|_1 [L_h(u, \mathcal{Z}_n)]^{1/m} du \\ &= N_h^{-D_b-1} \int_{U_n} \| N_h(t - \theta) - u \|_1 [L_h(\theta + u N_h^{-1}, \mathcal{Z}_n)]^{1/m} du \\ &= N_h^{-D_b-1} [L_h(\theta, \mathcal{Z}_n)]^{1/m} \int_{U_n} \| N_h(t - \theta) - u \|_1 Z_{n,\theta}^{1/m}(u) du. \end{aligned}$$

Hence,  $\tau_n = N_h(\hat{\theta}(h) - \theta)$  is the minimizer of

$$\chi_n(s) = \int_{U_n} \| s - u \|_1 z_n(u) du, \quad z_n(u) = \frac{Z_{n,\theta}^{1/m}(u)}{\int_{U_n} Z_{n,\theta}^{1/m}(v) dv},$$

and we obtain from (3.5.9) and (3.5.10) for any  $\varepsilon > 0$

$$(3.5.11) \quad \mathbb{P}_f \left\{ \| N_h(\hat{\theta}(h) - \theta) \|_1 > \varepsilon \right\} \leq \mathbb{P}_f \left\{ \inf_{\| s \|_1 > \varepsilon} \chi_n(s) \leq \chi_n(0) \right\}.$$

Let  $0 < r < \varepsilon/3$ , a number whose will be done later. We have

$$\chi_n(0) \leq r \int_{U_n \cap (\| u \|_1 \leq r)} z_n(u) du + \int_{U_n \cap (\| u \|_1 > r)} \| u \|_1 z_n(u) du.$$

Note also that

$$\begin{aligned} \inf_{\| s \|_1 > \varepsilon} \chi_n(s) &\geq \inf_{\| s \|_1 > \varepsilon} \left[ \int_{U_n \cap (\| u \|_1 \leq r)} (\| s \|_1 - \| u \|_1) z_n(u) du \right] \\ &\geq (\varepsilon - r) \int_{U_n \cap (\| u \|_1 \leq r)} z_n(u) du. \end{aligned}$$

It yields in particular

$$\chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) \leq -(\varepsilon - 2r) \int_{U_n \cap (\|u\|_1 \leq r)} z_n(u) du + \int_{U_n \cap (\|u\|_1 > r)} \|u\|_1 z_n(u) du.$$

Thus,  $\forall r \in (0, \varepsilon/3)$

$$\begin{aligned} & \mathbb{P}_f \left\{ \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0 \right\} \\ & \leq \mathbb{P}_f \left\{ \int_{U_n \cap (\|u\|_1 > r)} \|u\|_1 z_n(u) du > (\varepsilon - 2r) \int_{U_n \cap (\|u\|_1 \leq r)} z_n(u) du \right\} \\ & \leq \mathbb{P}_f \left\{ \int_{U_n \cap (\|u\|_1 > r)} \|u\|_1 z_n(u) du > r/2 \right\} \\ (3.5.12) \quad & + \mathbb{P}_f \left\{ (\varepsilon - 2r) \int_{U_n \cap (\|u\|_1 \leq r)} z_n(u) du < r/2 \right\}. \end{aligned}$$

We note that the second term in (3.5.12) can be controlled by the first one whenever  $0 < r < \varepsilon/3$ . Indeed, putting  $U_n(r) = U_n \cap (\|u\|_1 > r)$  we get

$$\begin{aligned} & \mathbb{P}_f \left\{ (\varepsilon - 2r) \int_{U_n \cap (\|u\|_1 \leq r)} z_n(u) du < r/2 \right\} \\ & \leq \mathbb{P}_f \left\{ r \int_{U_n} Z_{n,\theta}^{1/m}(v) dv - r \int_{U_n(r)} Z_{n,\theta}^{1/m}(u) du < \frac{r}{2} \int_{U_n} Z_{n,\theta}^{1/m}(v) dv \right\} \\ & \leq \mathbb{P}_f \left\{ r \int_{U_n(r)} Z_{n,\theta}^{1/m}(u) du > \frac{r}{2} \int_{U_n} Z_{n,\theta}^{1/m}(v) dv \right\} \\ & \leq \mathbb{P}_f \left\{ \int_{U_n(r)} \|u\|_1 z_n(u) du > r/2 \right\}. \end{aligned}$$

The last inequality and (3.5.12) yield

$$\mathbb{P}_f \left\{ \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0 \right\} \leq 2\mathbb{P}_f \left\{ \int_{U_n(r)} \|u\|_1 z_n(u) du > r/2 \right\}.$$

Let  $\nu > 0$  the parameter whose choice will be done later, then by definition of  $z_n$  and the last inequality, we obtain

$$\begin{aligned} & \mathbb{P}_f \left\{ \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0 \right\} \leq 2\mathbb{P}_f \left\{ \int_{U_n(r)} \|u\|_1 Z_{n,\theta}^{1/m}(u) du > r\nu/4 \right\} \\ (3.5.13) \quad & + 2\mathbb{P}_f \left\{ \int_{U_n} Z_{n,\theta}^{1/m}(v) dv < \nu/2 \right\}. \end{aligned}$$

In view of Lemma 7 and the Markov inequality, choosing  $r = \varepsilon/4 \geq s_2$  we get

$$(3.5.14) \quad \mathbb{P}_f \left\{ \int_{U_n \cap (\|u\|_1 > \varepsilon/4)} \|u\|_1 Z_{n,\theta}^{1/m}(u) du > \frac{\varepsilon\nu}{16} \right\} \leq 8C_2 I_1 \mathcal{E}_h^{c_2} \nu^{-1} e^{-c_3 2^{-2\gamma} \varepsilon^\gamma}.$$

Applying Lemma 6, we have for  $\nu \leq M$

$$(3.5.15) \quad \mathbb{P}_f \left\{ \int_{U_n} Z_{n,\theta}^{1/m}(v) dv < \nu/2 \right\} \leq 2C_1 \mathcal{E}_h^{c_1} \nu^{\frac{\tau}{D_b}}$$

In view of (3.5.13), (3.5.15), we establish that

$$\mathbb{P}_f \left\{ \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0 \right\} \leq 16C_2 I_1 \mathcal{E}_h^{c_2} \nu^{-1} e^{-c_3 2^{-2\gamma} \varepsilon^\gamma} + 4C_1 \mathcal{E}_h^{c_1} \mathcal{E}_h^{c_1} \nu^{\frac{\tau}{D_b}}.$$

Choosing  $\nu = e^{-c_3 2^{-2\gamma-1} \varepsilon^\gamma}$ , we obtain

$$\mathbb{P}_f \left\{ \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0 \right\} \leq 16C_2 I_1 \mathcal{E}_h^{c_2} e^{-c_3 2^{-2\gamma-1} \varepsilon^\gamma} + 4C_1 \mathcal{E}_h^{c_1} e^{-c_3 \frac{2^{-2\gamma-1} \tau}{D_b} \varepsilon^\gamma}.$$

The assertion of the theorem follows now from (3.5.9), (3.5.13) and the definitions of  $\omega_1$ ,  $\omega_2$ ,  $\omega_3$ . ■

### 3.6 Appendix

In this part, we consider  $nh^d \in \mathbb{N}^*$ , in order to simplify the proofs.

**Proof of Lemma 1.** 1. By definition of model 3.4.1, we have  $\forall u \in U_n$

$$Z_{n,\theta}(u) = \prod_{X_i \in V_h(y)} \exp \left\{ -\frac{1}{2} (Y_i - f_{\theta+u N_h^{-1}}(X_i))^2 + \frac{1}{2} (Y_i - f_\theta(X_i))^2 \right\}.$$

Let us give the notation

$$\Sigma = \sum_{X_i \in V_h(y)} \frac{1}{4} (Y_i - f_{\theta+u N_h^{-1}}(X_i))^2 - \frac{1}{4} (Y_i - f_\theta(X_i))^2.$$

Applying  $1 - e^{-\rho} \leq \rho$ ,  $\rho \in \mathbb{R}$  and Cauchy-Schwarz inequality, we obtain for any  $u : \|u\|_1 < 1$

$$\begin{aligned} \mathbb{E}_f (1 - Z_{n,\theta}^{1/2}(u))_+ &= \mathbb{E}_f (1 - e^{-\Sigma}) \mathbb{I}_{[0,+\infty[}(\Sigma) \\ &\leq \mathbb{E}_f(\Sigma) \mathbb{I}_{[0,+\infty[}(\Sigma) \\ &\leq 2\mathcal{E}_h \|u\|_1^2. \end{aligned}$$



Using the last inequality and define  $\Gamma_\delta = [0, \delta]^{D_b}$ , we establish

$$\int_{[0, \delta]^{D_b}} \mathbb{E}_f \left[ 1 - Z_{n, \theta}^{1/2}(u) \right]_+ du \leq 2\mathcal{E}_h \delta^{D_b+2}, \quad \forall \delta < 1.$$

2. Using  $\mathbb{E} e^{\lambda \varepsilon} = e^{\lambda^2/2}$ , with  $\varepsilon \sim \mathcal{N}(0, 1)$  and  $\lambda \in \mathbb{R}$ , we have

$$(3.6.1) \quad \mathbb{E}_f Z_{n, \theta}^{1/2}(u) = \exp \left\{ -\frac{1}{8} \sum_{X_i \in V_h(y)} f_{u N_h^{-1}}^2(X_i) + \frac{1}{2} \sum_{X_i \in V_h(y)} f_{u N_h^{-1}}(X_i) (f(X_i) - f_\theta(X_i)) \right\}.$$

Remark that  $f_{u N_h^{-1}}(X_i) \leq \|u\|_1 N_h^{-1}$ ,  $f(X_i) - f_\theta(X_i) \leq b_h$  and

$$\sum_{X_i \in V_h(y)} f_{u N_h^{-1}}^2(X_i) = u^\top \mathcal{M}_{nh}(y) u.$$

In view of Lemma 9 in Chapter 4,  $\exists \lambda > 0$  such that  $\lambda_n(h) \geq \lambda$  and  $u^\top \mathcal{M}_{nh}(y) u \geq \lambda_n(h) u^\top u$ . Applying (3.6.1), we have

$$\mathbb{E}_f Z_{n, \theta}^{1/2}(u) \leq \exp \left\{ -\frac{\lambda_n(h)}{8} \|u\|_2^2 + \frac{\|u\|_1}{2} \mathcal{N}(h) \right\},$$

By equivalence between the  $\ell_1$ -norm and  $\ell_2$ -norm, we obtain

$$\begin{aligned} \mathbb{E}_f Z_{n, \theta}^{1/2}(u) &\leq \exp \left\{ -\frac{D_b \lambda_n(h)}{8} \|u\|_1^2 + \frac{\|u\|_1}{2} \mathcal{N}(h) \right\} \\ &\leq \exp \left\{ -\frac{D_b \lambda_n(h)}{16} \|u\|_1^2 \right\}, \quad \forall u \in U_n : \|u\|_1 \geq \frac{8}{D_b \lambda_n(h)} \mathcal{N}(h). \end{aligned}$$

■

**Proof of Lemma 2.** 1. Remember that we have  $\forall u \in U_n = N_h(\Theta(A, M) - \theta)$

$$Z_{n, \theta}(u) = \prod_{X_i \in V_y(h)} \left( 1 + \frac{f_{u N_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^{Y_i} \exp \left\{ -f_{u N_h^{-1}}(X_i) \right\}, \quad N_h = (nh^d)^{1/2}.$$

Using the Cauchy-Schwarz inequality, we obtain for  $u \in U_n : \|u\|_1 < 1$

$$(3.6.2) \quad \mathbb{E}_f (1 - Z_{n, \theta}^{1/2}(u))_+ \leq \left( \mathbb{E}_f |1 - Z_{n, \theta}^{1/2}(u)|^2 \right)^{1/2} = \left( \mathbb{E}_f (1 + Z_{n, \theta}(u) - 2Z_{n, \theta}^{1/2}(u)) \right)^{1/2}.$$

It is easy to show that for  $n \geq n_0$

$$\mathbb{E}_f Z_{n, \theta}(u) \leq 1 + \frac{A + \mathcal{N}(h)}{A} \|u\|_1.$$

So, using (3.6.2) we have

$$(3.6.3) \quad \mathbb{E}_f(1 - Z_{n,\theta}^{1/2}(u))_+ \leq \left( \frac{A + \mathcal{N}(h)}{A} \|u\|_1 + 2(1 - \mathbb{E}_f Z_{n,\theta}^{1/2}(u)) \right)^{1/2}.$$

It remains to calculate the term

$$\begin{aligned} \mathbb{E}_f Z_{n,\theta}^{1/2}(u) &= \prod_{X_i \in V_h(y)} \mathbb{E}_f \left( 1 + \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^{Y_i/2} \exp \left\{ -f_{uN_h^{-1}}(X_i)/2 \right\} \\ &= \prod_{X_i \in V_h(y)} \sum_{k=0}^{\infty} \frac{(f(X_i))^k}{k!} \left( 1 + \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^{k/2} \exp \left\{ -f(X_i) - \frac{1}{2} f_{uN_h^{-1}}(X_i) \right\} \\ (3.6.4) \quad &= \prod_{X_i \in V_h(y)} \exp \left\{ -f(X_i) - \frac{1}{2} f_{uN_h^{-1}}(X_i) + f(X_i) \left( 1 + \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^{1/2} \right\} \end{aligned}$$

Using the Taylor expansion series, for  $n \geq n_0$  we get

$$\left( 1 + \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^{1/2} = 1 + \frac{1}{2} \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} - \frac{1}{4} \left( \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^2 + O \left[ \left( \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^3 \right],$$

where  $O(\rho^3) \xrightarrow{\rho \rightarrow 0} 0$ . Applying last equality and (3.6.4), we have

$$\begin{aligned} \mathbb{E}_f Z_{n,\theta}^{1/2}(u) &= \prod_{X_i \in V_h(y)} \exp \left\{ -\frac{1}{2} f_{uN_h^{-1}}(X_i) \frac{f_\theta(X_i) - f(X_i)}{f_\theta(X_i)} - \frac{1}{4} \left( f_{uN_h^{-1}}(X_i) \frac{f(X_i)}{f_\theta(X_i)} \right)^2 \right. \\ (3.6.5) \quad &\quad \left. + f(X_i) O \left[ \left( \frac{f_{uN_h^{-1}}(X_i)}{f_\theta(X_i)} \right)^3 \right] \right\}. \end{aligned}$$

Using the inequality  $1 - e^{-\rho} \leq \rho$ ,  $\rho \geq 0$ , (3.6.3), (3.6.4) and (3.6.5), we have for  $n \geq n_0$ ,

$$\mathbb{E}_f(1 - Z_{n,\theta}^{1/2}(u))_+ \leq \sqrt{\frac{2 + M + 2A}{2A}} \mathcal{E}_h^{1/2} \|u\|_1^{1/2}, \quad u \in [0, \delta]^{D_b}, \quad \delta < 1.$$

By integration, the first assertion of the lemma is proved.

**2.** Using (3.6.5) we have for  $n \geq n_0$  and  $\forall u \in U_n$  :  $\|u\|_1 \geq \frac{4M^2}{\lambda_n(h)} \mathcal{N}(h)$ , in view of Lemma 9 in Chapter 4,

$$\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq e^1 \exp \left\{ -\frac{\lambda_n(h)A}{8M^2} \|u\|_1^2 \right\},$$

where  $\lambda_n(h)$  is defined in Section 3.4. ■

**Proof of Lemma 3.** Remember that we have  $\forall u \in U_n = N_h(\Theta(M) - \theta)$

$$\begin{aligned} Z_{n,\theta}(u) &= \frac{L_h(\theta + u N_h^{-1}, \mathcal{Z}_n)}{L_h(\theta, \mathcal{Z}_n)} \\ &= \prod_{X_i \in V_y(h)} \exp \left\{ -|Y_i - f_{\theta + u N_h^{-1}}(X_i)|^\alpha + |Y_i - f_\theta(X_i)|^\alpha \right\}, \quad N_h = (nh^d)^{1/(1+2\alpha)}. \end{aligned}$$

It is easy to see that for  $m \geq 1$

$$\mathbb{E}_f Z_{n,\theta}^{1/m}(u) = \prod_{X_i \in V_h(y)} \int_{\mathbb{R}} C(\alpha) \exp \left\{ -\frac{1}{m}|y - b_i - v_i|^\alpha + \frac{1}{m}|y - b_i|^\alpha - |y|^\alpha \right\},$$

where  $b_i = f_\theta(X_i) - f(X_i)$  and  $v_i = f_{u N_h^{-1}}(X_i)$ . Remark that if we choose the approximation polynomial  $f_\theta$  such that  $f_\theta \geq f$  (for construction, see Chapter 4), then  $b_i \geq 0, \forall i$ . Without loss of generality we choose  $v_i \geq 0$ . Introduce

$$\mathcal{G}_m(v) = \int_{\mathbb{R}} C(\alpha) \exp \left\{ -\frac{1}{m}|y - b - v|^\alpha + \frac{1}{m}|y - b|^\alpha - |y|^\alpha \right\}, \quad v \geq 0, \quad b \geq 0.$$

We will prove that  $\mathcal{G}_m(v) \leq 1 - T_1 v^{1+2\alpha}$  for  $m > 1$ . Indeed, we use standard schema such that

$$(3.6.6) \quad \mathcal{G}_m(v) = \mathcal{G}_m(0) + \int_0^v \mathcal{G}'_m(\tilde{v}) d\tilde{v}.$$

It remains to prove that  $\mathcal{G}'_m(\tilde{v}) \leq -T_2 \tilde{v}^{2\alpha}$ . For it, we decompose the function  $\mathcal{G}_m$  of the following way

$$\begin{aligned} \mathcal{G}_m(v) &= \int_v^\infty C(\alpha) \exp \left\{ -\frac{1}{m}(y - v)^\alpha + \frac{1}{m}y^\alpha - (y + b)^\alpha \right\} dy \\ &\quad + \int_0^v C(\alpha) \exp \left\{ -\frac{1}{m}(v - y)^\alpha + \frac{1}{m}y^\alpha - (y + b)^\alpha \right\} dy \\ &\quad + \int_{-b}^0 C(\alpha) \exp \left\{ -\frac{1}{m}(v - y)^\alpha + \frac{1}{m}(-y)^\alpha - (y + b)^\alpha \right\} dy \\ &\quad + \int_{-\infty}^{-b} C(\alpha) \exp \left\{ -\frac{1}{m}(v - y)^\alpha + \frac{1}{m}(-y)^\alpha - (y + b)^\alpha \right\} dy \end{aligned}$$

The derivative of  $\mathcal{G}_m$  can be formulate like

$$\begin{aligned}
\mathcal{G}'_m(v) &= -C(\alpha) \exp \left\{ \frac{v^\alpha}{m} - (v+b)^\alpha \right\} + C(\alpha) \exp \left\{ \frac{v^\alpha}{m} - (v+b)^\alpha \right\} \\
&\quad + \int_v^\infty C(\alpha) \frac{\alpha}{m} (y-v)^{\alpha-1} \exp \left\{ -\frac{1}{m}(y-v)^\alpha + \frac{1}{m}y^\alpha - (y+b)^\alpha \right\} dy \\
&\quad - \int_0^v C(\alpha) \frac{\alpha}{m} (v-y)^{\alpha-1} \exp \left\{ -\frac{1}{m}(v-y)^\alpha + \frac{1}{m}y^\alpha - (y+b)^\alpha \right\} dy \\
&\quad - \int_{-b}^0 C(\alpha) \frac{\alpha}{m} (v-y)^{\alpha-1} \exp \left\{ -\frac{1}{m}(v-y)^\alpha + \frac{1}{m}(-y)^\alpha - (y+b)^\alpha \right\} dy \\
&\quad - \int_{-\infty}^{-b} C(\alpha) \frac{\alpha}{m} (v-y)^{\alpha-1} \exp \left\{ -\frac{1}{m}(v-y)^\alpha + \frac{1}{m}(-y)^\alpha - (y+b)^\alpha \right\} dy.
\end{aligned}$$

By changing variables and summation of integrals, we obtain

$$\begin{aligned}
(3.6.7) \quad \mathcal{G}'_m(v) &= \int_{v+b}^\infty C(\alpha) \frac{\alpha}{m} z^{\alpha-1} \exp \left\{ -\frac{1}{m}z^\alpha \right\} \left[ \exp \left\{ \frac{1}{m}(z+v)^\alpha - (z+v+b)^\alpha \right\} \right. \\
&\quad \left. - \exp \left\{ \frac{1}{m}(z-v)^\alpha - (z-v-b)^\alpha \right\} \right] dz
\end{aligned}$$

$$\begin{aligned}
(3.6.8) \quad &+ \int_v^{v+b} C(\alpha) \frac{\alpha}{m} z^{\alpha-1} \exp \left\{ -\frac{1}{m}z^\alpha \right\} \left[ \exp \left\{ \frac{1}{m}(z+v)^\alpha - (z+v+b)^\alpha \right\} \right. \\
&\quad \left. - \exp \left\{ \frac{1}{m}(z-v)^\alpha - (v-z+b)^\alpha \right\} \right] dz
\end{aligned}$$

$$\begin{aligned}
(3.6.9) \quad &+ \int_0^v C(\alpha) \frac{\alpha}{m} z^{\alpha-1} \exp \left\{ -\frac{1}{m}z^\alpha \right\} \left[ \exp \left\{ \frac{1}{m}(z+v)^\alpha - (z+v+b)^\alpha \right\} \right. \\
&\quad \left. - \exp \left\{ \frac{1}{m}(v-z)^\alpha - (v-z+b)^\alpha \right\} \right] dz.
\end{aligned}$$

By expansion Taylor series of order 1 and for any  $n \geq n_0$ , the term (3.6.9) can be controlled by

$$\begin{aligned}
&\int_0^v C(\alpha) \frac{\alpha}{m} z^{\alpha-1} \exp \left\{ -\frac{1}{m}z^\alpha \right\} \left( \frac{b^\alpha}{m} - \frac{m-1}{m}(v-z+b)^\alpha - \frac{m-1}{m}(z+v+b)^\alpha + O(b^{2\alpha}) \right) dz \\
&\leq (1 + C(\alpha)) \frac{v^\alpha b^\alpha}{m^2}.
\end{aligned}$$

The second term (3.6.8) can be bounded by  $(1 + C(\alpha)) \frac{v^\alpha b^\alpha}{m^2}$ . We note that the first term is the integral of a negative function, so we can increase it by

$$\begin{aligned} & \int_{v+b}^{2v+b} C(\alpha) \frac{\alpha}{m} z^{\alpha-1} \exp \left\{ -\frac{1}{m} z^\alpha \right\} \left( -\frac{m-1}{m} (z-v-b)^\alpha - \frac{m-1}{m} (z+v+b)^\alpha + O(v^{2\alpha}) \right) dz \\ & \leq - \int_{v+b}^{2v+b} C(\alpha) \frac{\alpha}{m} z^{\alpha-1} \exp \left\{ -\frac{1}{m} z^\alpha \right\} \frac{m-1}{m} (2z)^\alpha dz + O(v^{3\alpha}) \\ & \leq -v^{2\alpha} 2^{\alpha-2} C(\alpha) \frac{m-1}{m^2} \exp \left\{ -\frac{1}{m} v^\alpha \right\}. \end{aligned}$$

So we obtain

$$\mathcal{G}'_m(v) \leq 2(1 + C(\alpha)) \frac{v^\alpha b^\alpha}{m^2} - 2^{\alpha-2} C(\alpha) \frac{m-1}{m^2} \exp \left\{ -\frac{1}{m} v^\alpha \right\} v^{2\alpha},$$

Using (3.6.6), the last inequality implies that  $\forall m > 1$

$$(3.6.10) \quad \mathcal{G}_m(v) \leq 1 + 2(1 + C(\alpha)) \frac{b^\alpha |v|^{1+\alpha}}{(1+\alpha)m^2} - 2^{\alpha-2} C(\alpha) \frac{m-1}{(1+2\alpha)m^2} \exp \left\{ -\frac{1}{m} v^\alpha \right\} |v|^{1+2\alpha},$$

with  $\mathcal{G}_m(0) = 1$ . We can also show that it exists positive constants  $T_3, T_4$  such that

$$(3.6.11) \quad \mathcal{G}_1(v) \leq 1 + T_3 |v|^{1+2\alpha} + T_4 |v|^{1+\alpha} b^\alpha, \quad v \in \mathbb{R}.$$

1. Let us prove the first assumption,

$$\mathbb{E}_f (1 - Z_{n,\theta}^{1/2}(u))_+ \leq \mathbb{E}_f |1 - Z_{n,\theta}^{1/2}(u)| = \mathbb{E}_f \frac{1}{L_h^{1/2}(\theta, \mathcal{Z}_n)} \left| L_h^{1/2}(\theta, \mathcal{Z}_n) - L_h^{1/2}(\theta + u, \mathcal{Z}_n) \right|.$$

Let us note  $p_h(x, f) = \prod_{X_i \in V_h(y)} g_\alpha(x_i, f(X_i))$ , where  $x \in \mathbb{R}^{nh^d}$ . By integration and Cauchy-Schwarz inequality, we obtain that

$$\begin{aligned} \mathbb{E}_f (1 - Z_{n,\theta}^{1/2}(u))_+ & \leq \int_{\mathbb{R}^{nh^d}} \left| p_h^{1/2}(x, f_{\theta+u N_h^{-1}}) - p_h^{1/2}(x, f_\theta) \right| \frac{p_h(x, f)}{p_h^{1/2}(x, f_\theta)} dx \\ (3.6.12) \quad & \leq \left( \int_{\mathbb{R}^{nh^d}} \left| p_h^{1/2}(x, f_{\theta+u N_h^{-1}}) - p_h^{1/2}(x, f_\theta) \right|^2 dx \int_{\mathbb{R}^{nh^d}} \frac{p_h^2(x, f)}{p_h(x, f_\theta)} dx \right)^{1/2}. \end{aligned}$$

Introduce

$$H^2 \left( \mathbb{P}_{f_{\theta+u N_h^{-1}}}, \mathbb{P}_{f_\theta} \right) = \int_{\mathbb{R}^{nh^d}} \left| p_h^{1/2}(x, f_{\theta+u N_h^{-1}}) - p_h^{1/2}(x, f_\theta) \right|^2 dx,$$

where  $H(.,.)$  is the Hellinger distance between probability measures. More details are given in [Tsybakov \[2008\]](#) (Chapter 2). Remark that  $\mathbb{P}_{f_\theta}$  is a product measure  $\mathbb{P}_{f_\theta} = \otimes_{X_i \in V_h(y)} \mathbb{P}_{f_\theta}^i$  where  $\mathbb{P}_{f_\theta}^i = \mathbb{P}_{f_\theta(X_i)}$  the probability law of  $Y_i$ . By property of Hellinger distance, we obtain

$$(3.6.13) \quad H^2 \left( \mathbb{P}_{f_{\theta+u} N_h^{-1}}, \mathbb{P}_{f_\theta} \right) = 2 \left[ 1 - \prod_{X_i \in V_h(y)} \left( 1 - \frac{1}{2} H^2 \left( \mathbb{P}_{f_{\theta+u} N_h^{-1}}^i, \mathbb{P}_{f_\theta}^i \right) \right) \right],$$

where

$$\begin{aligned} H^2 \left( \mathbb{P}_{f_{\theta+u} N_h^{-1}}^i, \mathbb{P}_{f_\theta}^i \right) &= \int_{\mathbb{R}} \left| g_\alpha^{1/2}(x - f_{\theta+u} N_h^{-1}(X_i)) - g_\alpha^{1/2}(x - f_\theta(X_i)) \right|^2 dx, \\ &= \int_{\mathbb{R}} \left| g_\alpha^{1/2}(x - f_u N_h^{-1}(X_i)) - g_\alpha^{1/2}(x) \right|^2 dx, \quad i : X_i \in V_h(y). \end{aligned}$$

[Has'minskii and Ibragimov \[1981\]](#) showed that  $g_\alpha$  possess singularities of order  $2\alpha$  for  $\alpha < 1/2$ . In particular they gave the following result on  $g_\alpha$ :

$$(3.6.14) \quad \int_{\mathbb{R}} \left| g_\alpha^{1/2}(x - \eta) - g_\alpha^{1/2}(x) \right|^2 dx = O(\eta^{1+2\alpha}), \quad \eta \rightarrow 0.$$

For more details, see [Has'minskii and Ibragimov \[1981\]](#) (Chapter 6, section 1, Theorem 1.1). Using (3.6.12), (3.6.13) and 3.6.14, then it exists a constant  $\mathcal{C}$  such that for  $n \geq n_0$

$$(3.6.15) \quad \mathbb{E}_f (1 - Z_{n,\theta}^{1/2}(u))_+ \leq \mathcal{C} \|u\|_1^{\alpha+1/2} \left( \int_{\mathbb{R}^{nh^d}} \frac{p_h^2(x, f)}{p_h(x, f_\theta)} dx \right)^{1/2}.$$

In order to show the assumption 1.1, it remains to prove the following assertion

$$(3.6.16) \quad \left( \int_{\mathbb{R}^{nh^d}} \frac{p_h^2(x, f)}{p_h(x, f_\theta)} dx \right)^{1/2} \leq \exp \left\{ \frac{T_3 + T_4}{2} \mathcal{N}(h)^{1+2\alpha} \right\}.$$

For it, we note that

$$\int_{\mathbb{R}^{nh^d}} \frac{p_h^2(x, f)}{p_h(x, f_\theta)} dx = \prod_{X_i \in V_h(y)} \mathcal{G}_1(-b_i).$$

So, applying (3.6.11) with  $v = -b_i$ , we obtain (3.6.16). Using (3.6.15), the first assertion of the lemma is proved.

**2.** Using (3.6.10) with  $m = 2$  and  $v = u_i$ , we obtain

$$\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq \exp \left\{ \frac{1 + C(\alpha)}{2(1 + \alpha)} \mathcal{N}(h) - \frac{2^{\alpha-4} C(\alpha)}{(1 + 2\alpha) \exp \left\{ \frac{1}{2} (2M)^\alpha \right\}} \sum_{X_i \in V_h(y)} \frac{|f_u(X_i)|^{1+2\alpha}}{nh^d} \right\}.$$

Remark that  $\sum_{X_i \in V_h(y)} \frac{|f_u(X_i)|^{1+2\alpha}}{nh^d} \geq \sum_{X_i \in V_h(y)} \frac{|f_u(X_i)|^2}{nh^d \|u\|_1^{1-2\alpha}}$  and

$$\sum_{X_i \in V_h(y)} f_{u N_h^{-1}}^2(X_i) = u^\top \mathcal{M}_{nh}(y) u.$$

In view of Lemma 9 (in Chapter 4) and last inequality, we have

$$\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq \exp \left\{ \frac{1 + C(\alpha)}{2(1 + \alpha)} \mathcal{N}(h) \right\} \times \exp \left\{ - \frac{\lambda_n(h) 2^{\alpha-4} C(\alpha)}{(1 + 2\alpha) \exp \left\{ \frac{1}{2} (2M)^\alpha \right\}} \|u\|_1^{1+2\alpha} \right\},$$

where  $\lambda_n(h)$  is defined in Section 3.4. ■

**Proof of Lemma 5.** Note that for any  $k \geq \kappa + 1$  and by definition of  $\hat{k}$  (3.3.1)

$$\{\hat{k} = k\} = \cup_{l \geq k} \left\{ |\hat{f}^{(k-1)}(y) - \hat{f}^{(l)}(y)| > S_n(l) \right\}.$$

Note that  $S_n(l)$  is monotonically increasing in  $l$  and, therefore,

$$\begin{aligned} \{\hat{k} = k\} &\subseteq \left\{ |\hat{f}^{(k-1)}(y) - f(y)| > 2^{-1} S_n(k-1) \right\} \\ &\cup \left[ \cup_{l \geq k} \left\{ |\hat{f}^{(l)}(y) - f(y)| > 2^{-1} S_n(l) \right\} \right]. \end{aligned}$$

We come to the following inequality: for any  $k \geq \kappa + 1$

$$\begin{aligned} \mathbb{P}_f(\hat{k} = k) &\leq \mathbb{P}_f \left\{ |\hat{f}^{(k-1)}(y) - f(y)| > 2^{-1} S_n(k-1) \right\} \\ (3.6.17) \quad &+ \sum_{l \geq k} \mathbb{P}_f \left\{ |\hat{f}^{(l)}(y) - f(y)| > 2^{-1} S_n(l) \right\}. \end{aligned}$$

Note that the definition of  $S_n(l)$  yields

$$N_{h_l} S_n(l) = \left( s_2^\gamma \vee \frac{2^{2\gamma+1}(2\gamma + 2dq)}{c_3 \gamma (1 \wedge \tau D_b^{-1})} \right)^{\frac{1}{\gamma}} [1 + \ln(h_{\max}/h_l)]^{1/\gamma}.$$

Thus, applying Proposition 3 with  $J_1 = Ld\omega_2$  and  $\varepsilon = N_{h_l} S_n(l)$ , we obtain  $\forall l \geq k-1$

$$\begin{aligned} \mathbb{P}_f \left\{ |\hat{f}^{(l)}(y) - f(y)| > 2^{-1} S_n(l) \right\} &\leq \omega_1 \mathcal{E}_{h_l}^{\omega_2} [h_{\max}/h_l]^{-\frac{2\gamma+2dq}{\gamma}} \\ (3.6.18) \quad &= \omega_1 \mathcal{E}_{h_l}^{\omega_2} 2^{-l \frac{2\gamma+2dq}{\gamma}}. \end{aligned}$$

Note that  $b_{h_l} \leq Ld h_l^\beta$  and  $h_l \leq h_\kappa \leq h^*$  since  $f \in \mathbb{H}_d(\beta, L, M)$  and, therefore,

$$(3.6.19) \quad \mathcal{N}(h_l) \leq Ld n(h_l)^{\gamma\beta+d} \leq Ld n(h^*)^{\gamma\beta+d}, \quad \forall l \geq k-1.$$

Here we have also used that  $k \geq \kappa + 1$ . We obtain from (3.6.17), (3.6.18) and (3.6.19) that  $k \geq \kappa + 1$

$$\mathbb{P}_f(\hat{k} = k) \leq J_2 \exp \{J_1 n(h^*)^{\gamma\beta+d}\} 2^{-(k-1)\frac{2\gamma+2dq}{\gamma}},$$

where  $J_2 = \omega_1 \left(1 - 2^{-\frac{2\gamma+2dq}{\gamma}}\right)^{-1}$ . ■

**Proof of Lemma 6.** Applying Markov inequality, Fubini theorem and assumption 1.1

$$\begin{aligned} & \mathbb{P}_f \left\{ \int_{U_n} Z_{n,\theta}^{1/m}(u) du < \frac{1}{2} \delta^{D_b}, G \right\} \\ &= \mathbb{P}_f \left\{ \int_{\Gamma_\delta} (Z_{n,\theta}^{1/m}(u) - 1) du < -\frac{1}{2} \delta^{D_b}, \right\} \\ &\leq \mathbb{P}_f \left\{ \int_{\Gamma_\delta} (1 - Z_{n,\theta}^{1/m}(u))_+ du > \frac{1}{2} \delta^{D_b} \right\} \\ &\leq 2\delta^{-D_b} \int_{\Gamma_\delta} \mathbb{E}_f(1 - Z_{n,\theta}^{1/m}(u))_+ du \\ &\leq 2C_1 \mathcal{E}_h^{c_1} \delta^\tau. \end{aligned}$$
■

**Proof of Lemma 7.** Applying the assumption 1.2 we obtain

$$\mathbb{E}_f \int_{U_n \cap (\|u\|_1 > a)} \|u\|_1 Z_{n,\theta}^{1/m}(u) du \leq C_2 \mathcal{E}_h^{c_2} \int_{U_n \cap (\|u\|_1 > a)} \|u\|_1 \exp \{-c_3 \|u\|_1^\gamma\} du.$$

Simplifying the last inequality, we have

$$\mathbb{E}_f \{Q_a\} \leq 2C_2 a \mathcal{E}_h^{c_2} \exp \{-c_3 a^\gamma\} \int_0^{+\infty} (z+1) \exp \{-c_3 z^\gamma\} dz.$$
■





# Chapter 4

## Locally Bayesian Approach for Multiplicative Uniform Regression

Ce chapitre correspond à l'article [Chichignoud \[2010a\]](#). Dans cet article, nous présentons *l'estimateur Bayésien* sous une forme différente que celle dans le chapitre précédent. La démonstration des grandes déviations est plus complexe et longue que dans le chapitre précédent. Pour le lecteur ce chapitre contient la preuve des hypothèses [1](#) pour le modèle de régression avec un bruit multiplicatif uniforme (Voir lemme [8](#)) , ainsi que la borne inférieure du risque minimax (Section [4.2](#)) et l'optimalité de la vitesse adaptative (Section [4.3](#)) pour ce modèle.

### 4.1 Introduction

Let statistical experiment be generated by the couples of observations  $Y^{(n)} = (X_i, Y_i)_{i=1, \dots, n}, n \in \mathbb{N}^*$  where  $(X_i, Y_i)$  satisfies the equation

$$(4.1.1) \quad Y_i = f(X_i) \times U_i, \quad i = 1, \dots, n.$$

Here  $f : [0, 1]^d \rightarrow \mathbb{R}$  is unknown function and we are interested in estimating  $f$  at a given point  $y \in [0, 1]^d$  from observation  $Y^{(n)}$ .

The random variables (noise)  $(U_i)_{i=1, \dots, n}$  are supposed to be independent and uniformly distributed on  $[0, 1]$ .

The design points  $(X_i)_{i=1, \dots, n}$  are deterministic and without loss of generality we will assume that

$$X_i \in \{1/n^{1/d}, 2/n^{1/d}, \dots, 1\}^d, \quad i = 1, \dots, n.$$

Along the paper the unknown function  $f$  is supposed to be smooth, in particular, it belongs to the Holder ball of functions  $\mathbb{H}_d(\beta, L, M)$  (see Definition [13](#) below). Here  $\beta > 0$  is the smoothness of  $f$ ,  $M$  is the upper bound of  $f$  and its partial derivatives and  $L > 0$  is Lipschitz constant.

Moreover, we will consider only the functions  $f$  separated away from zero by some positive constant. Thus, from now on we will suppose that there exists  $0 < A < M$  such that  $f \in \mathbb{H}_d(\beta, L, M, A)$ , where

$$\mathbb{H}_d(\beta, L, M, A) = \left\{ g \in \mathbb{H}_d(\beta, L, M) : \inf_{x \in [0,1]^d} g(x) \geq A \right\}.$$

**Motivation.** The multiplicative regression model is quite popular in various domains of applications, in particular, in volatility estimation [Härlde and Tsybakov \[1997\]](#) or in noise speckle imaging [Aubert and Aujol \[2008\]](#). Another source of the interest to multiplicative regression is so-called *nonparametric frontier model* (see [Simar and Wilson \[2000\]](#)), where the reconstruction of the regression function  $f$  can be viewed as the estimation a production set  $\mathcal{P}$ . Indeed,  $Y_i \leq f(X_i)$ ,  $\forall i$ , and, therefore, the estimation of  $f$  is reduced to finding the upper boundary of  $\mathcal{P}$ . In this context one can also cite [Korostel'ev and Tsybakov \[1993\]](#) dealing with the estimation of function's support.

The theoretical interest to the multiplicative regression model (4.1.1) with discontinuous noise is dictated by the following fact. The typical approach to the study of the models with multiplicative noise consists in their transformation into the model with an additive noise and in the application, after that, the linear smoothing technique, based on standard methods like kernel smoothing, local polynomials etc. Let us illustrate the latter approach by the consideration of one of the most popular non-parametric model namely multiplicative gaussian regression

$$(4.1.2) \quad Y_i = \sigma(X_i)\xi_i, \quad i = 1, \dots, n.$$

Here  $\xi_i, i = 1, \dots, n$  are i.i.d. standard gaussian random variables and the goal is to estimate the variance  $\sigma^2(\cdot)$ .

Putting  $Y'_i = Y_i^2$  and  $\eta_i = \xi_i^2 - 1$  one can transform the model (4.1.2) into the heteroscedastic additive regression:

$$(4.1.3) \quad Y'_i = \sigma^2(X_i) + \sigma^2(X_i)\eta_i, \quad i = 1, \dots, n,$$

where, obviously,  $\mathbb{E}\eta_i = 0$ . Applying any of the linear methods mentioned above to the estimation of  $\sigma^2(\cdot)$  one can construct an estimator whose estimation accuracy is given by  $n^{-\frac{\beta}{2\beta+d}}$  and which is optimal in minimax sense (See Definition 14). The latter result is proved under assumptions on  $\sigma^2(\cdot)$  which are similar to the assumption imposed on the function  $f(\cdot)$ . In particular,  $\beta$  denotes the regularity of the function  $\sigma^2(\cdot)$ . The same result can be obtained for any noise variables  $\xi_i$  with known, continuously differentiable density, possessing sufficiently many moments.

The situation changes dramatically when one considers the noise with discontinuous distribution density. Although, the transformation of the original multiplicative model to the additive one is still possible, in particular, the model (4.1.1) can be rewritten as

$$Y'_i = f(X_i) + f(X_i)\eta_i, \quad Y'_i = 2Y_i, \quad \eta_i = 2u_i - 1, \quad i = 1, \dots, n,$$

the linear methods are not optimal anymore. As it is proved in Theorem 10 the optimal accuracy is given by  $n^{-\frac{\beta}{\beta+d}}$ . To achieve this rate the non-linear estimation procedure, based on locally bayesian approach, is proposed in Section 4.2.

The interesting feature of the model is that the uniform distribution leads to an improved estimation rate, it is always the case with a discontinuous density of the noise (see Has'minskii and Ibragimov [1981], Chapter 6). The parameter of a uniform  $[0, \theta]$  distribution is easier to estimate. Indeed, it is well known that we can use the maximum likelihood estimator based on the maximum of observations (nonlinear). For instance, if  $f$  is known to be Lipschitz  $\beta$ , then its approximation by a constant function in an hypercube of length  $h$  gives a bias of order  $dh^\beta$ , and the constant can be estimated at the order  $1/nh^d$ , because there are  $nh^d$  observations in the interval. The latter improved rate (improved over  $1/\sqrt{nh^d}$  as usual) leads to an optimal bandwidth of  $h = \bar{h}(\beta) = n^{-\frac{1}{\beta+d}}$ , and hence estimation rate  $n^{-\frac{\beta}{\beta+d}}$ , faster than usual. Typically, in any model, it is possible to improve the rate of convergence if the Fisher information is not finite.

In *regular statistical models*, the locally maximum likelihood approach is ideal (see Has'minskii and Ibragimov [1981] Chapter 1, Section 5 and recently Polzehl and Spokoiny [2006], Katkovnik and Spokoiny [2008]). But when the density of observations is discontinuous, an bayesian approach is often preferred (see Has'minskii and Ibragimov [1981], Chapter 1, Section 5). The study of local maximum likelihood estimator for model (4.1.1) stays an open problem.

Recently several approaches to the selection from the family of linear estimators were proposed, see for instance Goldenshluger and Lepski [2008], Goldenshluger and Lepski [2009a], Juditsky, Lepski, and Tsybakov [2009] and the references therein. These technologies are used in the construction of data-driven (adaptive) procedures and they are heavily based on the linearity property. As we already mentioned the locally bayesian estimators are non-linear and in Section 4.3 we propose the selection rule from this family. It requires, in particular, to develop new non-asymptotical exponential inequalities, which may have an independent interest.

**Minimax estimation.** The first part of the paper is devoted to the minimax over  $\mathbb{H}_d(\beta, L, M, A)$  estimation. This means, in particular, that the parameters  $\beta, L, M$  and  $A$  are supposed to be known *a priori*. We find the *minimax rate of convergence* (4.1.4) on  $\mathbb{H}_d(\beta, L, M, A)$  and propose the estimator being optimal in minimax sense (see Definition 14). Our first result in this direction consists in establishing a lower bound for maximal risk on  $\mathbb{H}_d(\beta, L, M, A)$ . We show that for any  $\beta \in \mathbb{R}_+^*$ , the minimax rate of convergence is bounded from below by the sequence

$$(4.1.4) \quad \varphi_n(\beta) = n^{-\frac{\beta}{\beta+d}}.$$

Next, we propose the minimax estimator, i.e. the estimator attaining the normalizing sequence (4.1.4). To construct the minimax estimator we use so-called *locally bayesian*

*estimation construction* which consists in the following. Let

$$V_h(y) = \bigotimes_{j=1}^d [y_j - h/2, y_j + h/2],$$

be the neighborhood around  $y$  such that  $V_h(y) \subseteq [0, 1]^d$ , where  $h \in (0, 1)$  is a given scalar. Fix an integer number  $b > 0$  and let

$$D_b = \sum_{m=0}^b \binom{m+d-1}{d-1}.$$

Let  $\mathcal{P}_b = \{p = (p_1, \dots, p_d) : p_i \in \mathbb{N}, 0 \leq |p| \leq b\}$ ,  $|p| = p_1 + \dots + p_d$ , we define the local polynomial

$$(4.1.5) \quad f_t(x) = \sum_{p \in \mathcal{P}_b} t_p \left( \frac{x-y}{h} \right)^p \mathbb{I}_{V_h(y)}(x), \quad x \in \mathbb{R}^d, t = (t_p : p \in \mathcal{P}_b),$$

where  $z^p = z_1^{p_1} \dots z_d^{p_d}$  for  $z = (z_1, \dots, z_d)$  and  $\mathbb{I}$  denotes the indicator function. The local polynomial  $f_t$  can be viewed as an approximation of the regression function  $f$  inside of the neighborhood  $V_h$  and  $D_b$  the number of coefficients of this polynomial. Introduce the following subset of  $\mathbb{R}^{D_b}$

$$(4.1.6) \quad \Theta(A, M) = \{t \in \mathbb{R}^{D_b} : 2t_{0,\dots,0} - \|t\|_1 \geq A, \|t\|_1 \leq M\},$$

where  $\|\cdot\|_1$  is  $l_1$ -norm on  $\mathbb{R}^{D_b}$ .  $\Theta(A, M)$  can be viewed as the set of coefficients  $t$  such that  $A \leq f_t(x) \leq M$  for all  $t \in \Theta(A, M)$  and for all  $x$  in the neighbourhood  $V_h(y)$ . Consider the *pseudo likelihood ratio*

$$L_h(t, Y^{(n)}) = \prod_{i: X_i \in V_h(y)} [f_t(X_i)]^{-1} \mathbb{I}_{[0, f_t(X_i)]}(Y_i), \quad t \in \Theta(A, M).$$

Set also

$$(4.1.7) \quad \pi_h(t) = \int_{\Theta(A, M)} \|t - u\|_1 L_h(u, Y^{(n)}) du, \quad t \in \Theta(A, M).$$

Let  $\hat{\theta}(h)$  be the solution of the following minimization problem:

$$(4.1.8) \quad \hat{\theta}(h) = \arg \min_{t \in \Theta(A, M)} \pi_h(t).$$

The *locally bayesian estimator*  $\bar{f}^h(y)$  of  $f(y)$  is defined now as  $\bar{f}^h(y) = \hat{\theta}_{0,\dots,0}(h)$ . Note that this local approach allows to estimate successive derivatives of function  $f$ . In this paper, only the estimation of  $f$  at a given point is studied.

We note that similar locally parametric approach based on maximum likelihood estimators was recently proposed in [Polzehl and Spokoiny \[2006\]](#) and [Katkovnik and Spokoiny \[2008\]](#) for *regular statistical models*.

As we see our construction contains an extra-parameter  $h$  to be chosen. To make this choice we use quite standard arguments. First, we note that in view of  $f \in \mathbb{H}_d(\beta, L, M, A)$

$$\exists \theta = \theta(f, y, h) \in \Theta(A, M) : \sup_{x \in V_h(y)} |f(x) - f_\theta(x)| \leq Lh^\beta.$$

For example,  $f_\theta$  can be chosen like the Taylor polynomial (defined in (4.5.2)). With the property of isotropic Holder functions (see Definition 13), the last assertion is trivial. Thus, if  $h$  is chosen sufficiently small our original model (4.1.1) is well approximated inside of  $V_h(y)$  by the "parametric" model

$$\mathcal{Y}_i = f_\theta(X_i) \times U_i, \quad i = 1, \dots, nh^d, \quad nh^d \in \mathbb{N}^*$$

in which the *bayesian estimator*  $\hat{\theta}$  is rate-optimal (See Theorem 11).

It is worth mentioning that the analysis of the deviation of  $(X_i, \mathcal{Y}_i)_{i=1, \dots, nh^d}$  from  $Y^{(nh^d)}$  is not simple. Namely here requirements  $0 < A \leq f(x) \leq M, \forall x \in [0, 1]^d$ , are used. This assumption, which seems not to be necessary, allows us to make the presentation of basic ideas clear and to simplify routine computations (see also Remark 12). Without this assumption, the problem is open.

Finally,  $h_n(\beta, L)$  is chosen as the solution of the following minimization problem

$$(4.1.9) \quad (nh^d)^{-1} + Lh^\beta \rightarrow \min_h$$

and we show that corresponding estimator  $\bar{f}^{h_n(\beta, L)}(y)$  is minimax for  $f(y)$  on  $\mathbb{H}_d(\beta, L, M, A)$  if  $\beta \leq b$ . Since the parameter  $b > 0$  can be chosen in arbitrary way, the proposed estimator is minimax for any given value of the parameter  $\beta > 0$ .

**Adaptive estimation.** The second part of the paper is devoted to the adaptive minimax estimation over collection of isotropic functional classes in the model (4.1.1). At our knowledge, the problem of adaptive estimation in the multiplicative regression with the noise, having discontinuous density, is not studied in the literature.

Well-known drawback of minimax approach is the dependence of the minimax estimator on the parameters describing functional class on which the maximal risk is determined. In particular, the locally bayesian estimator  $\bar{f}^h(\cdot)$  depends obviously on the parameters  $A$  and  $M$  via the solution of the minimization problem (4.1.8). Moreover  $h_n(\beta, L)$  optimally chosen in view of (4.1.9) depends explicitly on  $\beta$  and  $L$ . To overcome this drawback the minimax adaptive approach was proposed (see [Lepski \[1990\]](#), [Lepski \[1991\]](#), [Lepski, Mammen, and Spokoiny \[1997\]](#)). The first question arising in the adaptation (reduced to the problem at hand) can be formulated as follows.

*Does there exist an estimator which would be minimax on  $\mathbb{H}(\beta, L, M, A)$  simultaneously for all values of  $\beta, L, A$  and  $M$  belonging to some given subset of  $\mathbb{R}_+^4$  ?*

In section 4.3, we show that the answer to this question is **negative**, that is typical for the estimation of the function at a given point [Lepski and Spokoiny \[1997\]](#). This answer can be reformulated in the following manner: the family of rates of convergence  $\{\varphi_n(\beta), \beta \in \mathbb{R}_+^*\}$  is **unattainable** for the problem under consideration.

Thus, we need to find another family of normalizations for maximal risk which would be attainable and, moreover, optimal in view of some criterion of optimality. Nowadays, the most developed criterion of optimality is due to *Klutchnikoff* [Klutchnikoff \[2005\]](#).

We show that the family of normalizations, being optimal in view of this criterion, is

$$(4.1.10) \quad \phi_n(\beta) = \left( \frac{\rho_n(\beta)}{n} \right)^{\frac{\beta}{\beta+d}}, \quad \rho_n(\beta) = 1 + \ln \left( \frac{\varphi_n(\beta)}{\varphi_n(b)} \right),$$

whenever  $\beta \in ]0, b]$ . The factor  $\rho_n$  can be considered as *price to pay for adaptation* [Lepski \[1990\]](#).

The most important step in proving the optimality of the family (4.1.10) is to find an estimator, called *adaptive*, which attains the optimal family of normalizations. Obviously, we seek an estimator whose construction is *parameter-free*, i.e. independent of  $\beta, L, A$  and  $M$ . In order to explain our estimation procedure let us make several remarks.

First we note that the role of the constants  $A, M$  and  $\beta, L$  in the construction of the minimax estimator is quite different. Indeed, the constants  $A, M$  are used in order to determine the set  $\Theta(A, M)$  needed for the construction of the locally bayesian estimator, see (4.1.7) and (4.1.8). However, this set does not depend on the localization parameter  $h > 0$ , in other words, the quantities  $A$  and  $M$  are not involved in the selection of optimal size of the local neighborhood given by (4.1.9). Contrary to that, the constants  $\beta, L$  are used for the derivation of the optimal size of the local neighborhood (4.1.9), but they are not involved in the construction of the collection of locally bayesian estimators  $\{\hat{f}^h, h > 0\}$ .

Next remark explains how to replace the unknown quantities  $A$  and  $M$  in the definition of  $\Theta(A, M)$ . Our first simple observation consists in the following: the estimator  $\bar{f}^{h_n(\beta, L)}$  remains minimax if we replace  $\Theta(A, M)$  in (4.1.7) and (4.1.8) by  $\Theta(\tilde{A}, \tilde{M})$  with any  $0 < \tilde{A} \leq A$  and  $M \leq \tilde{M} < \infty$ . It follows from obvious inclusion  $\mathbb{H}_d(\beta, L, A, M) \subseteq \mathbb{H}_d(\beta, L, \tilde{A}, \tilde{M})$ . The next observation is less trivial and it follows from Proposition 4. Put  $h_{\max} = n^{-\frac{1}{b+d}}$  and define for any function  $f$

$$(4.1.11) \quad A(f) = \inf_{x \in V_{h_{\max}}(y)} f(x), \quad M(f) = \sum_{m=0}^b \sum_{p_1+\dots+p_d=m} \left| \frac{\partial^m f(y)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} \right|.$$

The following agreement will be used in the sequel: if the function  $f$  and  $m \geq 1$  be such that  $\partial^m f$  does not exist we will put formally  $\partial^m f = 0$  in the definition of  $M(f)$ .

Then,  $\bar{f}^{h_n(\beta, L)}$  remains minimax if we replace  $\Theta(A, M)$  in (4.1.7) by  $\Theta(A(f), M(f))$ . It remains to note that contrary to the quantities  $A$  and  $M$  the functionals  $A(f)$  and  $M(f)$  can be consistently estimated from the observation (4.1.1) and let  $\hat{A}$  and  $\hat{M}$  be the corresponding estimators. The idea now is to determine the collection of locally bayesian estimators  $\{\hat{f}^h, h > 0\}$  by replacing  $\Theta(A, M)$  in (4.1.7) and (4.1.8) by the *random* parameter set  $\hat{\Theta}$  which is defined as follows.

$$\hat{\Theta} = \left\{ t \in \mathbb{R}^{D_b} : 2t_{0, \dots, 0} - \|t\|_1 \geq 2^{-1}\hat{A}, \quad \|t\|_1 \leq 4\hat{M} \right\}.$$

In this context it is important to emphasize that the estimators  $\hat{A}$  and  $\hat{M}$  are built from the same observation which is used for the construction of the family  $\{\hat{f}^h, h > 0\}$ .

Contrary to all saying above, the constants  $\beta$  and  $L$  cannot be estimated consistently. In order to select an "optimal" estimator from the family  $\{\hat{f}^h, h > 0\}$  we use general adaptation scheme due to *Lepski Lepski* [1990], *Lepski* [1992a]. To the best of our knowledge it is the first time when this method is applied in the statistical model with multiplicative noise and discontinuous distribution. Moreover, except already mentioned papers *Polzehl and Spokoiny* [2006] and *Katkovnik and Spokoiny* [2008], Lepski's procedure is typically applied to the selection from the collection of linear estimators (kernel estimators, locally polynomial estimator, etc.). In the present paper we apply this method to very complicated family of nonlinear estimators, obtained by the use of bayesian approach on the random parameter set. It required, in particular, to establish the exponential inequality for the deviation of locally bayesian estimator from the parameter to be estimated (Proposition 4). It generalizes the inequality proved for the parametric model *Has'minskii and Ibragimov* [1981] (Chapter 1, Section 5). This result seems to be new.

**Simulations.** The interest to the multiplicative models is explained by their popularity in the image processing, for example, SAR images *Réfrégier* [2002]. If the model is parametric then the use of the maximum likelihood estimator is rather standard at least if the likelihood is convex *Bhattachar and Sundareshan* [2000]. In the present paper we adopt the local parametric approximation to a purely non parametric model. As it proved, this strategy leads to the theoretically optimal statistical decisions. But the minimax as well as the minimax adaptive approach are asymptotical and it seems natural to check how proposed estimators work for reasonable sample size. In the simulation study, we test the bayesian estimator in the parametric and nonparametric cases. We show that the *adaptive* estimator approaches the *oracle* estimator. The *oracle* estimator is selected from the family  $\{\hat{f}^h, h > 0\}$  under the hypothesis  $f$  that is known. We show that the bayesian estimator performs well starting with  $n \geq 100$ .

This paper is organized as follows. In Section 4.2 we present the results concerning minimax estimation and Section 4.3 is devoted to the adaptive estimation. The simulations



are given in Section 4.4. The proofs of main results are proved in Section 4.5 (upper bounds) and section 4.6 (lower bounds). Auxiliary lemmas are postponed to Appendix (Section 4.7) contains the proofs of technical results.

## 4.2 Minimax estimation on isotropic Hölder class

In this section we present several results concerning minimax estimation. First, we establish lower bound for minimax risk defined on  $\mathbb{H}_d(\beta, L, M, A)$  for any  $\beta, L, M$  and  $A$ . For any  $(p_1, \dots, p_d) \in \mathbb{N}^d$  we denote  $p = (p_1, \dots, p_d)$  and  $|p| = p_1 + \dots + p_d$ .

**Definition 13.** Fix  $\beta > 0$ ,  $L > 0$  and  $M > 0$  and let  $\lfloor \beta \rfloor$  be the largest integer strictly less than  $\beta$ . The isotropic Hölder class  $\mathbb{H}_d(\beta, L, M)$  is the set of functions  $f : [0, 1]^d \rightarrow \mathbb{R}$  having on  $[0, 1]^d$  all partial derivatives of order  $\lfloor \beta \rfloor$  and such that  $\forall x, y \in [0, 1]^d$

$$\sum_{m=0}^{\lfloor \beta \rfloor} \sum_{|p|=m} \sup_{x \in [0, 1]^d} \left| \frac{\partial^{|p|} f(x)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} \right| \leq M,$$

$$\left| \frac{\partial^{|p|} f(x)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} - \frac{\partial^{|p|} f(y)}{\partial y_1^{p_1} \dots \partial y_d^{p_d}} \right| \leq L [\|x - y\|_1]^{\beta - \lfloor \beta \rfloor}, \quad \forall |p| = \lfloor \beta \rfloor.$$

This definition implies that if  $f \in \mathbb{H}_d(\beta, L, M, A)$ , then  $A \leq A(f)$  and  $M(f) \leq M$ , where  $\mathbb{H}_d(\beta, L, M, A)$  is defined in the beginning of this paper.  $A(f)$  and  $M(f)$  are defined in (4.1.11).

**Maximal and minimax risk on  $\mathbb{H}_d(\beta, L, M, A)$ .** To measure the performance of estimation procedures on  $\mathbb{H}_d(\beta, L, M, A)$  we will use minimax approach.

Let  $\mathbb{E}_f = \mathbb{E}_f^n$  be the mathematical expectation with respect to the probability law of the observation  $Y^{(n)}$  satisfying (4.1.1). We define first the maximal risk on  $\mathbb{H}_d(\beta, L, M, A)$  corresponding to the estimation of the function  $f$  at a given point  $y \in [0, 1]^d$ .

Let  $\tilde{f}$  be an arbitrary estimator built from the observation  $Y^{(n)}$ . Let  $\forall q > 0$

$$R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M, A)] = \sup_{f \in \mathbb{H}_d(\beta, L, M, A)} \mathbb{E}_f |\tilde{f}(y) - f(y)|^q.$$

The quantity  $R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M, A)]$  is called *maximal risk* of the estimator  $\tilde{f}$  on  $\mathbb{H}_d(\beta, L, M, A)$  and the *minimax risk* on  $\mathbb{H}_d(\beta, L, M, A)$  is defined as

$$R_{n,q}[\mathbb{H}_d(\beta, L, M, A)] = \inf_{\tilde{f}} R_{n,q}[\tilde{f}, \mathbb{H}_d(\beta, L, M, A)],$$

where inf is taken over the set of all estimators.

**Definition 14.** The normalizing sequence  $\psi_n$  is called minimax rate of convergence (MRT) and the estimator  $\hat{f}$  is called minimax (asymptotically minimax) if

$$\begin{aligned} \liminf_{n \rightarrow \infty} \psi_n^{-q} R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M, A)] &> 0; \\ \limsup_{n \rightarrow \infty} \psi_n^{-q} R_{n,q}[\hat{f}, \mathbb{H}_d(\beta, L, M, A)] &< \infty. \end{aligned}$$

**Theorem 10.** For any  $\beta > 0$ ,  $L > 0$ ,  $M > 0$ ,  $A > 0$ ,  $q \geq 1$  and  $d \geq 1$

$$\liminf_{n \rightarrow \infty} \varphi_n^{-q}(\beta) R_{n,q}[\mathbb{H}_d(\beta, L, M, A)] > 0, \quad \varphi_n(\beta) = n^{-\frac{\beta}{\beta+d}}.$$

**Remark 12.** The obtained result shows that on  $\mathbb{H}_d(\beta, L, M, A)$  the minimax rate of convergence cannot be faster than  $n^{-\frac{\beta}{\beta+d}}$ . In view of the obvious inclusion  $\mathbb{H}_d(\beta, L, M, A) \subset \mathbb{H}_d(\beta, L, M)$  the minimax rate of convergence on an isotropic Hölder class is also bounded from below by  $n^{-\frac{\beta}{\beta+d}}$ .

The next theorem shows how to construct the minimax estimator basing on locally bayesian approach. Put  $\bar{h} = (Ln)^{-\frac{1}{\beta+d}}$  and let  $\bar{f}^{\bar{h}}(y) = \hat{\theta}_{0,\dots,0}(\bar{h})$  is given by (4.1.6), (4.1.7) and (4.1.8) with  $h = \bar{h}$ .

**Theorem 11.** Let  $\beta > 0$ ,  $L > 0$ ,  $M > 0$  and  $A > 0$  be fixed. Then there exists the constant  $C_*$  such that for any  $n \in \mathbb{N}^*$  satisfying  $n\bar{h}^d \geq (\lfloor \beta \rfloor + 1)^d$

$$\varphi_n^{-q}(\beta) R_{n,q}[\bar{f}^{\bar{h}}(y), \mathbb{H}_d(\beta, L, M, A)] \leq C^*, \quad \forall q \geq 1.$$

The explicit form of  $C^*$  is given in the proof.

**Remark 13.** We deduce from Theorems 10 and 11 that the estimator  $\bar{f}^{\bar{h}}(y)$  is minimax on  $\mathbb{H}_d(\beta, L, M, A)$ .

### 4.3 Adaptive estimation on isotropic Hölder classes

This section is devoted to the adaptive estimation over the collection of the classes  $\left\{ \mathbb{H}_d(\beta, L, M, A) \right\}_{\beta, L, M, A}$ . We will not impose any restriction on possible values of  $L, M, A$ , but we will assume that  $\beta \in (0, b]$ , where  $b$ , as previously, is an arbitrary a priori chosen integer.

We start with formulating the result showing that there is no optimally adaptive estimator (here we follow the terminology introduced in Lepski [1991], Lepski [1992a]). It means that there is no an estimator which would be minimax simultaneously for several values of parameter  $\beta$  even if all other parameters  $L, M$  and  $A$  are supposed to be fixed. This result does not require any restriction on  $\beta$  as well.

**Theorem 12.** For any subset of  $\mathbb{R}^+ \setminus \{0\}$  noted  $\mathbb{B}$  such that  $\text{card}(\mathbb{B}) \geq 2$ , for any  $\beta_1, \beta_2 \in \mathbb{B}$  and any  $L > 0, M > 0, A > 0$

$$\liminf_{n \rightarrow \infty} \inf_{\tilde{f}} \left[ \varphi_n^{-q}(\beta_1) R_{n,q}(\tilde{f}, \mathbb{H}_d(\beta_1, L, M, A)) + \varphi_n^{-q}(\beta_2) R_{n,q}(\tilde{f}, \mathbb{H}_d(\beta_2, L, M, A)) \right] = +\infty,$$

where  $\inf$  is taken over all possible estimators.

The assertion of Theorem 12 can be considerably specified if  $\mathbb{B} = (0, b]$ . To do that we will need the following definition. Let  $\Psi = \{\psi_n(\beta)\}_{\beta \in (0, b]}$  be a given family of normalizations.

**Definition 15.** The family  $\Psi$  is called admissible if there exist an estimator  $\hat{f}_n$  such that for some  $L > 0, M > 0$  and  $A > 0$

$$(4.3.1) \quad \limsup_{n \rightarrow \infty} \psi_n^{-q}(\beta) R_{n,q}(\hat{f}_n, \mathbb{H}_d(\beta, L, M, A)) < \infty, \quad \forall \beta \in (0, b].$$

The estimator  $\hat{f}_n$  satisfying (4.3.1) is called  $\Psi$ -attainable. The estimator  $\hat{f}_n$  is called  $\Psi$ -adaptive if (4.3.1) holds for any  $L > 0, M > 0$  and  $A > 0$ .

Note that the result proved in Theorem 12 means that the family of rates of convergence  $\{\varphi_n(\beta)\}_{\beta \in (0, b]}$  is not admissible. Let  $\Phi$  be the following family of normalizations:

$$\phi_n(\beta) = \left( \frac{\rho_n(\beta)}{n} \right)^{\frac{\beta}{\beta+d}}, \quad \rho_n(\beta) = 1 + \ln \left( \frac{\varphi_n(\beta)}{\varphi_n(b)} \right), \quad \beta \in (0, b].$$

We remark that  $\phi_n(b) = \varphi_n(b)$  and  $\rho_n(\beta) \sim \ln n$  for any  $\beta \neq b$ .

**Theorem 13.** Let  $\Psi = \{\psi_n(\beta)\}_{\beta \in (0, b]}$  be an arbitrary admissible family of normalizations.

**I.** For any  $\alpha \in (0, b]$  such that  $\psi_n(\alpha) \neq \varphi_n(\alpha)$ , there exists an admissible family  $\{v_n(\beta)\}_{\beta \in (0, b]}$  for which

$$\lim_{n \rightarrow \infty} v_n(\alpha) \psi_n^{-1}(\alpha) = 0.$$

**II.** If there exists  $\gamma \in (0, b)$  such that

$$(4.3.2) \quad \lim_{n \rightarrow \infty} \psi_n(\gamma) \phi_n^{-1}(\gamma) = 0,$$

then necessarily

$$(a) \quad \lim_{n \rightarrow \infty} \psi_n(\beta) \phi_n^{-1}(\beta) > 0, \quad \forall \beta \in (0, \gamma);$$

$$(b) \quad \lim_{n \rightarrow \infty} \left[ \frac{\psi_n(\gamma)}{\phi_n(\gamma)} \right] \left[ \frac{\phi_n(\beta)}{\psi_n(\beta)} \right] = 0, \quad \forall \beta \in (\gamma, b].$$

Several remarks are in order.

1. We note that if the family of normalizations  $\Phi$  is admissible, i.e. one can construct  $\Phi$ -attainable estimator, then  $\Phi$  is in an *optimal* family of normalizations in view of Kluchnikoff criterion Klutchnikoff [2005]. It follows from the second assertion of the theorem. We note however that a  $\Phi$ -attainable estimator may depend on  $L > 0$ ,  $M > 0$  and  $A > 0$ , and, therefore, this estimator have only theoretical interest. In the next section we construct  $\Phi$ -adaptive estimator, which is, by its definition, fully parameter-free. Moreover, this estimator obviously proves that  $\Phi$  is admissible, and, therefore, optimal as it was mentioned above.
2. The assertions of Theorem 13 allows us to give rather simple interpretation of Kluchnikoff criterion. Indeed, the first assertion, which is easily deduced from Theorem 12, shows that any admissible family of normalizations can be improved by another admissible family at any given point  $\alpha \in (0, b]$  except maybe one. In particular, it concerns the family  $\Phi$  if it is admissible. On the other hand, the second assertion of the theorem shows that there is no admissible family which would outperform the family  $\Phi$  at two points. Moreover, in view of (b),  $\Phi$ -adaptive (attainable) estimator, if exists, has the same precision on  $\mathbb{H}_d(\beta, L, M, A)$ ,  $\beta < \gamma$ , as any  $\Psi$ -adaptive(attainable) estimator whenever  $\Psi$  satisfies (4.3.2). Additionally, (a) implies that the gain in the precision provided by  $\Psi$ -adaptive (attainable) estimator on  $\mathbb{H}_d(\gamma, L, M, A)$  leads automatically to much more losses on  $\mathbb{H}_d(\beta, L, M, A)$  for any  $\beta > \gamma$  with respect to the precision provided by  $\Phi$ -adaptive(attainable) estimator. We conclude that  $\Phi$ -adaptive(attainable) estimator outperforms any  $\Psi$ -adaptive(attainable) estimator whenever  $\Psi$  satisfies (4.3.2). It remains to note that any admissible family not satisfying (4.3.2) is asymptotically equivalent to  $\Phi$ .

**Construction of  $\Phi$ -adaptive estimator.** As it was already mentioned in Introduction the construction of our estimation procedure consists of several steps. First, we determine the set  $\hat{\Theta}$ , built from observation, which is used after that in order to define the family of locally bayesian estimators. Next, based on Lepski's method (see Lepski [1991] and Lepski, Mammen, and Spokoiny [1997]), we propose data-driven selection from this family.

First step: Determination of parameter set. Put  $h_{\max} = n^{-\frac{1}{b+d}}$  and let  $\tilde{\theta}$  be the solution of the following minimization problem.

$$(4.3.3) \quad \inf_{t \in \mathbb{R}^{D_b}} \sum_{i: X_i \in V_{\max}(y)}^n \left[ 2Y_i - t K^\top \left( \frac{X_i - y}{h_{\max}} \right) \right]^2, \quad V_{\max}(y) = V_{h_{\max}}(y),$$

where the  $D_b$ -dimensional vector  $K(z) = (z^p : p \in \mathcal{P}_b)$  and the sign  $\top$  below means the transposition. Thus,  $\tilde{\theta}$  is the local least squared estimator and its explicit expression is given

by

$$(4.3.4) \quad \tilde{\theta} = 2 \left[ \sum_{i: X_i \in V_{\max}(y)}^n K^\top \left( \frac{X_i - y}{h_{\max}} \right) K \left( \frac{X_i - y}{h_{\max}} \right) \right]^{-1} [\mathcal{K}_n(y)]^\top Y,$$

where  $Y = (Y_1, \dots, Y_n)$  and  $\mathcal{K}_n(y) = \left[ K^\top \left( \frac{X_i - y}{h_{\max}} \right) \mathbb{I}_{V_{\max}(y)}(X_i) \right]_{i=1, \dots, n}$  is the design matrix. Put

$$\tilde{\delta}_p = p_1! \dots p_d! h_{\max}^{-|p|} \tilde{\theta}_p, \quad |p| \leq b.$$

Introduce the following quantities

$$(4.3.5) \quad \hat{A} = \tilde{\delta}_{0, \dots, 0}, \quad \hat{M} = \|\tilde{\delta}\|_1,$$

and define the random parameter set as follows.

$$(4.3.6) \quad \hat{\Theta} = \left\{ t \in \mathbb{R}^{D_b} : 2t_{0, \dots, 0} - \|t\|_1 \geq 2^{-1} \hat{A}, \quad \|t\|_1 \leq 4\hat{M} \right\}.$$

Second step: Collection of locally bayesian estimators. Put

$$(4.3.7) \quad \hat{\pi}_h(t) = \int_{\hat{\Theta}} \|t - u\|_1 L_h(u, Y^{(n)}) du;$$

$$(4.3.8) \quad \hat{\theta}^*(h) = \arg \min_{t \in \hat{\Theta}} \hat{\pi}_h(t).$$

The family of locally bayesian estimator  $\hat{\mathcal{F}}$  is defined now as follows.

$$(4.3.9) \quad \hat{\mathcal{F}} = \left\{ \hat{f}^h(y) = \hat{\theta}_{0, \dots, 0}^*(h), \quad h \in (0, h_{\max}] \right\}.$$

Third step: Data-driven selection from the collection  $\hat{\mathcal{F}}$ . Put

$$h_k = 2^{-k} h_{\max}, \quad k = 0, \dots, k_n,$$

where  $k_n$  is smallest integer such that  $h_{k_n} \geq h_{\min} = \ln^{\frac{b}{d(b+d)}} n^{-1/d}$ . Set

$$\hat{\mathcal{F}}^* = \left\{ \hat{f}^{(k)}(y) = \hat{\theta}_{0, \dots, 0}^*(h_k), \quad k = 0, \dots, k_n \right\}.$$

We put  $\hat{f}^*(y) = \hat{f}^{(\hat{k})}(y)$ , where  $\hat{f}^{(\hat{k})}(y)$  is selected from  $\hat{\mathcal{F}}^*$  in accordance with the rule:

$$(4.3.10) \quad \hat{k} = \inf \left\{ k = \overline{0, k_n} : |\hat{f}^{(k)}(y) - \hat{f}^{(l)}(y)| \leq \hat{M} S_n(l), \quad l = \overline{k+1, k_n} \right\}.$$

Here we have used the following notations.

$$(4.3.11) \quad S_n(l) = 432 D_b^3 (32qd + 16) \lambda_n^{-1}(h_l) \left[ \frac{1 + l \ln 2}{n(h_l)^d} \right], \quad l = 0, 1, \dots, k_n,$$

and  $\lambda_n(h)$  is the smallest eigenvalue of the matrix

$$(4.3.12) \quad \mathcal{M}_{nh}(y) = \frac{1}{nh^d} \sum_{i=1}^n K^\top \left( \frac{X_i - y}{h} \right) K \left( \frac{X_i - y}{h} \right) \mathbb{I}_{V_h(y)}(X_i),$$

which is completely determined by the design points and by the number of observations. We will prove that there exists a nonnegative real  $\lambda$ , such that  $\lambda_n(h) \geq \lambda$  for any  $n \geq 1$  and any  $h \in [h_{\min}, h_{\max}]$  (see Lemma 9).

**Theorem 14.** *Let an integer number  $b > 0$  be fixed. Then for any  $\beta \in (0, b]$ ,  $L > 0$ ,  $M > 0$ ,  $A > 0$  and  $q \geq 1$*

$$\limsup_{n \rightarrow \infty} \phi_n^{-q}(\beta) R_{n,q} \left[ \hat{f}^*(y), \mathbb{H}_d(\beta, L, M, A) \right] < \infty.$$

**Remark 14.** *The assertion of the theorem means that the proposed estimator  $\hat{f}^*(y)$  is  $\Phi$ -adaptive. It implies in particular that the family of normalizations  $\Phi$  is admissible. This, together with Theorem 13 allows us to state the optimality of  $\Phi$  in view of Kluchnikoff criterion (see Klutchnikoff [2005]).*

## 4.4 Simulation study

We will consider the case  $d = 1$ . The data are simulated accordingly to the model (4.1.1), where we use the following functions (Figure 4.1).

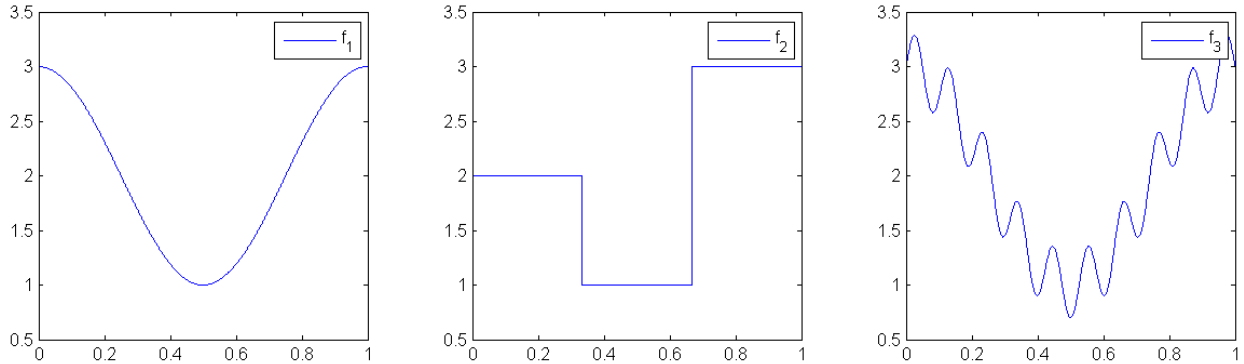


Figure 4.1:

Here  $f_1(x) = \cos(2\pi x) + 2$ ,  $f_2(x) = 2.\mathbb{I}_{[x \leq 1/3]} + 1.\mathbb{I}_{[1/3 < x \leq 2/3]} + 3.\mathbb{I}_{[2/3 < x]}$  and  $f_3(x) = \cos(2\pi x) + 2 + 0.3 \sin(19\pi x)$

To construct the family of estimators we use the linear approximation ( $b = 2$ ), i.e. within the neighbourhoods of the given size  $h$ , the locally bayesian estimator has the form

$$\hat{f}^h(x) = \hat{\theta}_0 + \hat{\theta}_1 x, \quad x \in [0, 1].$$

We define the ideal (oracle) value of the parameter  $\tilde{h} = \tilde{h}(f)$  as the minimizer of the risk:

$$\tilde{h} = \arg \inf_{h \in [1/n, 1]} \mathbb{E}_f |\hat{f}^h(y) - f(y)|.$$

To compute it we apply Monte-Carlo simulations (10000 repetitions). Our first objective is to compare the risk provided by the "oracle" estimator  $\hat{f}^{\tilde{h}}(\cdot)$  and whose provided by the adaptive estimator from Section 4.3. Figure 4.2 shows the deviation of the adaptive estimator from the function to be estimated. In several points, for example in  $y = 1/2$ , we remark so-called over-smoothing phenomenon, inherent to any adaptive estimator.

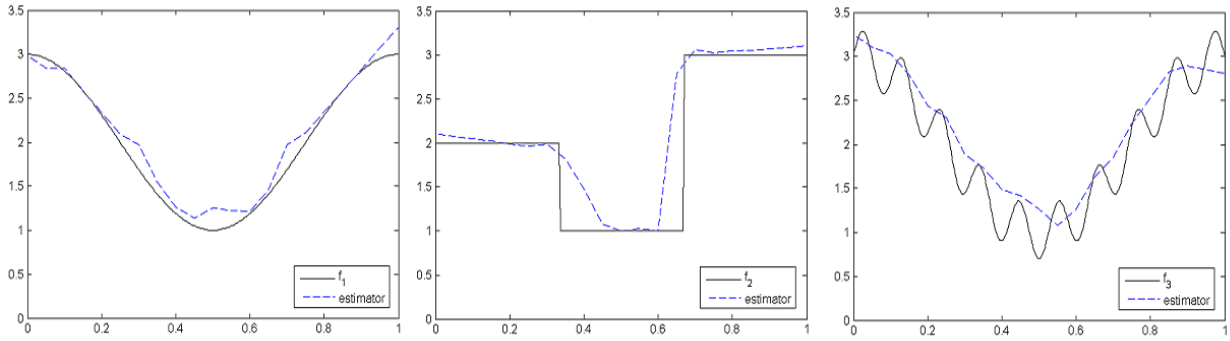


Figure 4.2:

**Oracle-adaptive ratio.** We compute the risks of the oracle and the adaptive estimator in 100 points of the interval  $(0, 1)$ . The next tabular presents the mean value of the ratio oracle risk/adaptive risk calculated for the functions  $f_1, f_2, f_3$  and  $n = 100, 1000$ .

function	n = 100		n = 1000	
	adaptive risk	oracle-adaptive ratio	adaptive risk	oracle-adaptive ratio
$f_1$	0.13	0.84	0.03	0.85
$f_2$	0.3	0.71	0.1	0.75
$f_3$	0.28	0.65	0.2	0.68

Figure 4.3 presents the "oracle risk/adaptive risk" ratio as the function of the number of observations  $n$ .

**Adaptation versus parametric estimation.** We consider the function  $f_4$  (figure 4.4), which is linear inside the neighborhood of size  $h_* = 1/8$  around point  $1/2$  and simulate  $n = 1000$  observations in accordance with the model (4.1.1). Using only the observations corresponding to the interval  $[3/8, 5/8]$  we construct the bayesian estimator  $\hat{f}^{1/8}(1/2)$ .

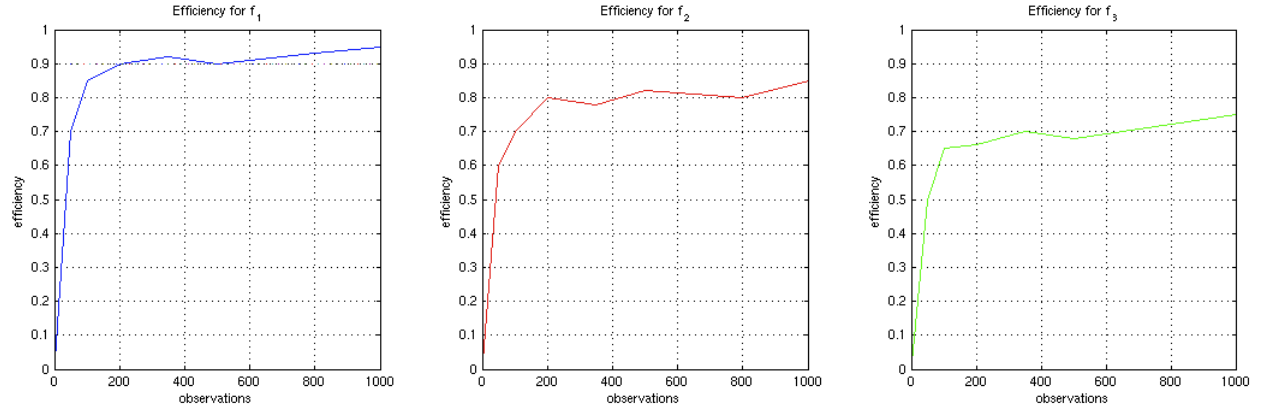


Figure 4.3:

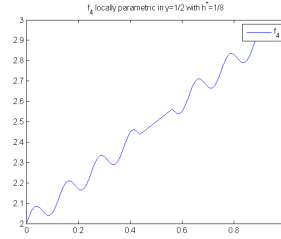


Figure 4.4:

It is important to emphasize that this estimator is efficient [Has'minskii and Ibragimov \[1981\]](#) since the model is parametric. Our objective now is to compare the risk of our adaptive estimator with the risk provided by the estimator  $\hat{f}^{1/8}(1/2)$ . We also try to understand how far is the localization parameter  $h_{\hat{k}}$ , inherent to the construction of our adaptive estimator, from the true value  $1/8$ . We compute the risk of each estimator via Monte-Carlo method with 10000 repetitions. For each repetition the procedure select the adaptive bandwidth  $h_{\hat{k}}^{(j)}$ ,  $j = 1, \dots, 10000$ . We confirm once again the over-smoothing phenomenon since

$$h_{\hat{k}}^{(j)} \sim 0.1405 > h_* = 0.1250, \quad j = 1, \dots, 10000.$$

Note however that the adaptive procedure selects the neighborhood of the size which is quite close to the true one. We also compute the risks of both estimators: "bayesian risk"=**0.0206** and "adaptive risk"=**0.0308**". We conclude that the estimation accuracy provided by our adaptive procedure is quite satisfactory.



## 4.5 Proofs of main results: upper bounds

Let  $\mathcal{H}_n, n > 1$  be the following subinterval of  $(0, 1)$ .

$$(4.5.1) \quad \mathcal{H}_n = \left[ \frac{(b+1) \vee (\ln n)^{\frac{1}{(d+d^2)}}}{n^{1/d}}, \left( \frac{1}{\ln n} \right)^{\frac{1}{b+d}} \right].$$

Later on we will consider only the values of  $h$  belonging to  $\mathcal{H}_n$ . We begin by stating the control of large deviations which is necessary to prove minimax and adaptive results.

### 4.5.1 Auxiliary results

Introduce the following notations. For any  $h > 0$  satisfying (4.5.1) put  $\omega = \omega(f, y, h) = \{\omega_p : p \in \mathcal{P}_b\}$ , where  $\omega_0 = \omega_{0,\dots,0} = f(y)$  and

$$(4.5.2) \quad \omega_p = \frac{\partial^{|p|} f(y)}{\partial y_1^{p_1} \cdots \partial y_d^{p_d}} \frac{h^{|p|}}{p_1! \cdots p_d!}, \quad p \in \mathcal{P}_b.$$

Remind the agreement which we follow in the present paper: if the function  $f$  and vector  $p$  are such that  $\partial^{|p|} f$  does not exist we put  $\omega_p = 0$ .

Let  $f_\omega(x)$ , given by (4.1.5), be the local polynomial approximation of  $f$  inside  $V_h(y)$  and let  $b_h$  be the corresponding approximation error, i.e.

$$(4.5.3) \quad b_h = \sup_{x \in V_h(y)} |f_\omega(x) - f(x)|.$$

It is easy to show that  $b_h \leq Ldh^\beta$  by definition of  $\omega$  in (4.5.2). Put finally the bias-variance ratio and an other term which are necessary for the interpretation of the following results.

$$(4.5.4) \quad \mathcal{N}_h = b_h \times nh^d, \quad \mathcal{E}(h) = \exp \left\{ \frac{(1 + 6D_b^2)\mathcal{N}_h}{6A(f)D_b^2} \right\}.$$

Introduce the random events  $G_{\hat{M}} = \{|\hat{M} - M(f)| \leq M(f)/2\}$  and  $G_{\hat{A}} = \{|\hat{A} - A(f)| \leq A(f)/2\}$  and put  $G = G_{\hat{M}} \cap G_{\hat{A}}$  where  $\hat{A}$  and  $\hat{M}$  are defined in (4.3.5) in Section 4.3. Recall that  $\lambda_n(h)$  (defined in Section 4.3) is the smallest eigenvalue of the matrix

$$\mathcal{M}_{nh}(y) = \frac{1}{nh^d} \sum_{i=1}^n K^\top \left( \frac{X_i - y}{h} \right) K \left( \frac{X_i - y}{h} \right) \mathbb{I}_{V_h(y)}(X_i),$$

and  $K(z)$  is the  $D_b$ -dimensional vector of the monomials  $z^p$ ,  $p \in \mathcal{P}_b$ .

The next proposition is the milestone for all results proved in the paper.

**Proposition 4.** *For any  $h$  satisfying (4.5.1) and any  $f$  such that  $A(f) > 0$  and  $M(f) < \infty$ , then  $\forall \varepsilon > 0$*

$$\mathbb{P}_f \left( nh^d |\hat{f}^h(y) - f(y)| \geq \varepsilon, G \right) \leq \mathfrak{B}(A(f), M(f)) \mathcal{E}(h) \exp \left\{ -\frac{\lambda_n(h) \varepsilon}{432 M(f) D_b^3} \right\},$$

where  $\hat{f}^h(y) \in \hat{\mathcal{F}}$  defined in (4.3.9). The explicit expression of the function  $\mathfrak{B}(\cdot, \cdot)$  is given in the beginning of the proof of the proposition.

In seeking this type of result (control of large deviations), the new point is to control the term  $\mathcal{E}(h)$  which can be very large compared to the normalization  $(nh^d)^{-1}$ . It is a consequence that the bayesian estimator is non-linear. The bias-variance ratio  $\mathcal{N}_h$  (defined in (4.5.4)), which appears in  $\mathcal{E}(h)$ , allows to obtain similar results with linear estimators (including the choice of the bandwidth).

The next proposition provides us with upper bound for the risk of a locally bayesian estimator.

**Proposition 5.** *For any  $n \in \mathbb{N}^*$ ,  $h \in \mathcal{H}_n$  and any  $f \in \mathbb{H}_d(\beta, L, M, A)$ , then  $\exists \lambda > 0$  does not depend of  $n$  such that:*

$$\mathbb{E}_f |\hat{f}^h(y) - f(y)|^q \mathbb{I}_G \leq C_q^* \left[ \frac{1 \vee L d n h^{\beta+d}}{n h^d} \right]^q, \quad q \geq 1.$$

where

$$C_q^* = \frac{\mathfrak{B}(A, M)}{q} \left( \frac{144 M(f) D_b}{A(f) \lambda (1 + 6 D_b^2)^{-1}} \right)^q \int_0^{+\infty} (\eta + 1)^{q-1} \exp \left\{ -\frac{\lambda \eta}{432 M(f) D_b^3} \right\} d\eta.$$

### 4.5.2 Proof of Proposition 4

Before beginning this proof, we introduce several notations which simplify understanding and we state Lemma 8. The knowledge of the likelihood  $L_h$  is not directly involved in the proof of Proposition 1, but only in the proof of Lemma 8. In fact, the proof of Proposition 4 requires only Lemma 8 and using the bayesian estimator with the local non-parametric approach. For it, we use standard arguments like Markov inequality or the following inclusion

$$\left\{ nh^d \|\hat{\theta}^*(h) - \theta\|_1 \geq \varepsilon \right\} \subseteq \left\{ \inf_{nh^d \|t - \theta\|_1 \geq \varepsilon} \hat{\pi}_h(t) \leq \hat{\pi}_h(\theta) \right\}.$$

where  $\hat{\theta}^*(h)$  minimizes  $\hat{\pi}_h$  defined in (4.3.7).

Set for  $z > 0, u > 0$  and  $v > 0$  Define finally for any  $u > 0, v > 0$

$$(4.5.5) \quad \mathfrak{B}(u, v) = \sup_{z \geq 0} 16e (1 \vee D_b u^{-1}) \Sigma^*(v) [\mathcal{B}_z + 6] \exp \left\{ -\frac{\lambda z}{432 v D_b^3} \right\},$$

where  $\mathcal{B}_z = z^{D_b+1} + 2^{D_b}(2z+2)^{\frac{D_b}{2}} + 6 + D_b \left( z^{D_b} + (2z+2)^{\frac{D_b}{2}-1} \right)$ ,  $\lambda > 0$  is described in Lemma 9 and does not depend of  $n$ , and

$$(4.5.6) \quad \Sigma^*(v) = \frac{c^2(v)(3-c(v))}{(1-c(v))^3}, \quad c(v) = \exp \{ -(54vD_b^2)^{-1} \}.$$

**Auxiliary Lemma.** Introduce the vector  $\theta = \theta(f, y, h) = \{\theta_p : p \in \mathcal{P}_b\}$ , where

$$(4.5.7) \quad \theta_0 = \theta_{0,\dots,0} = \omega_0 + b_h, \quad \theta_p = \omega_p, \quad |p| \neq 0,$$

where  $\omega$  is the coefficients of Taylor polynomial defined in (4.5.2). The definition of  $b_h$  implies obviously

$$(4.5.8) \quad f_\theta(x) \geq f(x), \quad \forall x \in V_h(y).$$

This trivial inequality will be extensively exploited in the sequel.

Let  $U_n = \{u \in \mathbb{R}^{D_b} : u = nh^d(t - \theta), t \in \Theta(A(f)/4, 9M(f))\}$ . For any  $u \in U_n$ , put the process

$$Z_{h,\theta}(u) = \frac{L_h(\theta + (nh^d)^{-1}u, Y^{(n)})}{L_h(\theta, Y^{(n)})}.$$

Note that in view of (4.5.8), the event  $Y_i \leq f_\theta(X_i)$  is always realized, because the multiplicative noise belongs to  $[0, 1]$ , therefore  $Y_i \leq f(X_i) \leq f_\theta(X_i)$ . By definition of  $L_h$ , the process  $Z_{h,\theta}$  can be written as follows

$$(4.5.9) \quad Z_{h,\theta}(u) = \prod_{i: X_i \in V_h(y)} \frac{f_\theta(X_i)}{f_{\theta+u(nh^d)^{-1}}(X_i)} \mathbb{I}_{[Y_i \leq f_{\theta+u(nh^d)^{-1}}(X_i)]}, \quad u \in U_n.$$

**Lemma 8.** For any  $f \in \mathbb{H}_d(\beta, L, M, A)$  and  $h \in \mathcal{H}_n$

1.  $\sup_{u_1, u_2 \in U_n} \|u_1 - u_2\|_1^{-1} \mathbb{E}_f |Z_{h,\theta}(u_1) - Z_{h,\theta}(u_2)| \leq \mathcal{C}_h,$
2.  $\mathbb{E}_f Z_{n,\theta}^{1/2}(u) \leq e^{-g_h(\|u\|_1)}, \quad \forall u \in U_n,$
3.  $\mathbb{P}_f \left\{ \int_{[0,\delta]^{D_b}} Z_{h,\theta}(u) du < \frac{\delta^{D_b}}{2} \right\} < 2\mathcal{C}_h\delta, \quad \forall \delta > 0.$

where

$$\mathcal{C}_h = 8(1 \vee D_b A^{-1}(f)) \exp \{1 + \mathcal{N}_h/A(f)\}, \quad g_h(a) = \frac{\lambda_n(h)a}{18M(f)D_b} - \frac{\mathcal{N}_h}{A(f)}.$$

$\lambda_n(h)$  is the smallest eigenvalue of the matrix  $\mathcal{M}_{nh}(y)$  defined in (4.3.12).

The knowledge of the likelihood occurs only in the proof of previous lemma.

**Proof of Proposition 4.** This proof is based on a result, already known in nonparametric estimation, given by [Has'minskii and Ibragimov \[1981\]](#) in their book (Chapter 1, Section 5, Theorem 5.2). We begin by stating two assertions which will prove in the first and second parts of this proof. Finally, we conclude easily with these two assertions and the Markov inequality. Let us introduce a new process  $z_h$  which is the renormalized process  $Z_{h,\theta}$ .

$$z_h(u) = \frac{Z_{h,\theta}(u)}{\int_{\hat{U}_n} Z_{h,\theta}(v) dv}, \quad u \in \hat{U}_n$$

where  $\hat{U}_n = nh^d(\hat{\Theta} - \theta)$  where the set  $\hat{\Theta}$  is defined in (4.3.6). Remind that  $\theta = \theta(f, y, h)$  is defined in (4.5.7).

**Assertion 1.** For any  $\varepsilon > 0$ , and for all  $r$  such that  $0 < r < \varepsilon/3$ , we assume

$$\mathbb{P}_f \left( nh^d |\hat{f}^h(y) - f(y)| \geq \varepsilon, G \right) \leq 2\mathbb{P}_f \left( \int_{\hat{U}_n(r)} \|u\|_1 z_h(u) du > \frac{r}{2}, G \right).$$

The proof of Assertion 1 use the definition of the process  $z_h$  and the following inclusion

$$\left\{ nh^d \|\hat{\theta}^*(h) - \theta\|_1 \geq \varepsilon \right\} \subseteq \left\{ \inf_{nh^d \|t - \theta\|_1 \geq \varepsilon} \hat{\pi}_h(t) \leq \hat{\pi}_h(\theta) \right\}.$$

where  $\hat{\theta}^*(h)$  minimizes  $\hat{\pi}_h$  defined in (4.3.7).

**Assertion 2.** For all  $h \in \mathcal{H}_n$  and any  $f$  such that  $A(f) > 0$  and  $M(f) < \infty$ , then for any  $a > 0$

$$\mathbb{E}_f \left[ \int_{\hat{U}_n \cap \{\|u\|_1 > a\}} \|u\|_1 z_h(u) du \mathbb{I}_G \right] \leq a \Sigma^*(M(f)) \mathcal{B}_a C_h \exp \left\{ -\frac{1}{6D_b} g_h(a) \right\},$$

where  $g_h(\cdot)$  is defined in Lemma 8.

This assertion requires Markov inequality and Lemma 8.

*First step: Proof of Assertion 1.*

The definition of  $\hat{\theta}^*(h)$  and  $\theta = \theta(f, y, h)$  implies  $\forall \varepsilon > 0$

$$\begin{aligned} \mathbb{P}_f \left( nh^d |\hat{f}^h(y) - f(y)| \geq \varepsilon, G \right) &\leq \mathbb{P}_f \left( nh^d |\hat{\theta}_0^*(h) - \theta_0| \geq \varepsilon, G \right) \\ (4.5.10) \quad &\leq \mathbb{P}_f \left( nh^d \|\hat{\theta}^*(h) - \theta\|_1 \geq \varepsilon, G \right). \end{aligned}$$

Some remarks are in order. First, it is easily seen that  $\theta \in \Theta(A(f), 3M(f))$ . Therefore, if the event  $G$  holds then  $\theta \in \hat{\Theta}$ . Remind also that  $\hat{\theta}^*(h)$  minimizes  $\hat{\pi}_h$  defined in (4.3.7) and, therefore, the following inclusion holds since  $\hat{\theta}^*(h) \in \hat{\Theta}$ .

$$(4.5.11) \quad \left\{ \left( nh^d \|\hat{\theta}^*(h) - \theta\|_1 \geq \varepsilon \right) \cap G \right\} \subseteq \left\{ \left( \inf_{nh^d \|t - \theta\|_1 \geq \varepsilon} \hat{\pi}_h(t) \leq \hat{\pi}_h(\theta) \right) \cap G \right\}.$$

Moreover,

$$\begin{aligned}
\hat{\pi}_h(t) &= (nh^d)^{-1} \int_{\hat{\Theta}} \|nh^d(t-u)\|_1 L_h(u, Y^{(n)}) du \\
&= (nh^d)^{-D_b-1} \int_{\hat{U}_n} \|nh^d(t-\theta) - u\|_1 L_h(\theta + u(nh^d)^{-1}, Y^{(n)}) du \\
&= (nh^d)^{-D_b-1} L_h(\theta, Y^{(n)}) \int_{\hat{U}_n} \|nh^d(t-\theta) - u\|_1 Z_{h,\theta}(u) du.
\end{aligned}$$

Hence,  $\tau_n = nh^d(\hat{\theta}^*(h) - \theta)$  is the minimizer of

$$\chi_n(s) = \int_{\hat{U}_n} \|s - u\|_1 \frac{Z_{h,\theta}(u)}{\int_{\hat{U}_n} Z_{h,\theta}(v) dv} du$$

and we obtain from (4.5.10) and (4.5.11) for any  $\varepsilon > 0$

$$(4.5.12) \quad \mathbb{P}_f \left( \|nh^d(\hat{\theta}^*(h) - \theta)\|_1 > \varepsilon, G \right) \leq \mathbb{P}_f \left( \inf_{\|s\|_1 > \varepsilon} \chi_n(s) \leq \chi_n(0), G \right).$$

Let  $0 < r < \varepsilon/3$ , be a number whose choice will be done later. We have

$$\chi_n(0) \leq r \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} z_h(u) du + \int_{\hat{U}_n \cap (\|u\|_1 > r)} \|u\|_1 z_h(u) du.$$

Note also that

$$\begin{aligned}
\inf_{\|s\|_1 > \varepsilon} \chi_n(s) &\geq \inf_{\|s\|_1 > \varepsilon} \left[ \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} (\|s\|_1 - \|u\|_1) z_h(u) du \right] \\
&\geq (\varepsilon - r) \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} z_h(u) du.
\end{aligned}$$

It yields in particular

$$\begin{aligned}
\chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) &\leq -(\varepsilon - 2r) \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} z_h(u) du + \int_{\hat{U}_n \cap (\|u\|_1 > r)} \|u\|_1 z_h(u) du.
\end{aligned}$$

Thus,  $\forall r \in (0, \varepsilon/3)$

$$\begin{aligned}
&\mathbb{P}_f \left( \chi_n(0) - \inf_{\|s\|_1 > \varepsilon} \chi_n(s) > 0, G \right) \\
&\leq \mathbb{P}_f \left( \int_{\hat{U}_n \cap (\|u\|_1 > r)} \|u\|_1 z_h(u) du > (\varepsilon - 2r) \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} z_h(u) du, G \right) \\
&\leq \mathbb{P}_f \left( \int_{\hat{U}_n \cap (\|u\|_1 > r)} \|u\|_1 z_h(u) du > r/2, G \right) \\
(4.5.13) \quad &+ \mathbb{P}_f \left( (\varepsilon - 2r) \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} z_h(u) du < r/2, G \right).
\end{aligned}$$

We note that the second term in (4.5.13) can be control by the first one whenever  $0 < r < \varepsilon/3$ . Indeed, putting  $\hat{U}_n(r) = \hat{U}_n \cap (u \in \mathbb{R}^{D_b} : \|u\|_1 > r)$  we get

$$\begin{aligned} & \mathbb{P}_f \left( (\varepsilon - 2r) \int_{\hat{U}_n \cap (\|u\|_1 \leq r)} z_h(u) du < r/2, G \right) \\ & \leq \mathbb{P}_f \left( r \int_{\hat{U}_n} Z_{h,\theta}(v) dv - r \int_{\hat{U}_n(r)} Z_{h,\theta}(u) du < \frac{r}{2} \int_{\hat{U}_n} Z_{h,\theta}(v) dv, G \right) \\ & \leq \mathbb{P}_f \left( r \int_{\hat{U}_n(r)} Z_{h,\theta}(u) du > \frac{r}{2} \int_{\hat{U}_n} Z_{h,\theta}(v) dv, G \right) \\ & \leq \mathbb{P}_f \left( \int_{\hat{U}_n(r)} \|u\|_1 z_h(u) du > r/2, G \right). \end{aligned}$$

The last inequality together with (4.5.10), (4.5.12) and (4.5.13) yields

$$(4.5.14) \quad \mathbb{P}_f \left( nh^d |\hat{f}^h(y) - f(y)| \geq \varepsilon, G \right) \leq 2\mathbb{P}_f \left( \int_{\hat{U}_n(r)} \|u\|_1 z_h(u) du > \frac{r}{2}, G \right).$$

*Second step: Proof of Assertion 2.*

Put  $\overline{U_n(a)} = U_n \cap (u \in \mathbb{R}^{D_b} : \|u\|_1 > a)$  for all  $a > 0$  and  $\Gamma_v = U_n(v+1) \setminus U_n(v)$  for any  $v > a$ . Introduce the following notations.

$$\mathcal{I}_v = \int_{\Gamma_v} Z_{h,\theta}(u) du, \quad \mathcal{Q}_v = \frac{\mathcal{I}_v}{\int_{\hat{U}_n} Z_{h,\theta}(u) du}.$$

We divide the set  $\Gamma_v$  into  $T^{D_b}$  nonintersecting parts  $\Delta_1, \Delta_2, \dots, \Delta_{T^{D_b}}$  such that

$$\Gamma_v = \bigcup_{j=1}^{T^{D_b}} \Delta_j, \quad \Delta_j \cap \Delta_l = \emptyset, \quad \forall j, l.$$

The integer number  $T$  will be chosen later. Let  $u_i$  be the point arbitrary chosen in the set  $\Delta_i$ . Then

$$S_v = \sum_i \int_{\Delta_i} Z_{h,\theta}(u_i) du.$$

In fact, we seek to approach the process  $\mathcal{I}_v$  with a process in steps  $S_v$ . Obviously  $\bar{\Gamma}_v = \int_{\Gamma_v} du$  and we get for any  $\sigma > 0$

$$\begin{aligned} \mathbb{P}_f(S_v > \sigma) & \leq \mathbb{P}_f \left( \max_i Z_{h,\theta}^{1/2}(u_i) \sqrt{\bar{\Gamma}_v} > \sqrt{\sigma} \right) \\ & \leq \sum_i \mathbb{P}_f \left( Z_{h,\theta}^{1/2}(u_i) > (\bar{\Gamma}_v)^{-1/2} \sqrt{\sigma} \right). \end{aligned}$$

Note that the number of summands on the right-hand side of the last inequality is equal to  $T^{D_b}$ . Applying Markov inequality and Lemma 8.2, we obtain

$$(4.5.15) \quad \mathbb{P}_f(S_v > \sigma) \leq T^{D_b} \sqrt{\Gamma_v} \sigma^{-1/2} e^{-g_h(v)}.$$

In view of Lemma 8.1,

$$\mathbb{E}_f |S_v - \mathcal{I}_v| \leq \sum_i \int_{\Delta_i} \mathbb{E}_f |Z_{h,\theta}(u) - Z_{h,\theta}(u_i)| du \leq C_h \sum_i \int_{\Delta_i} \|u - u_i\|_1 du.$$

Note that each summand does not exceed  $(vT^{-1})^{D_b+1}$  and, therefore,

$$(4.5.16) \quad \mathbb{E}_f |S_v - \mathcal{I}_v| \leq C_h v^{D_b+1} T^{-1}.$$

One has

$$\mathbb{P}_f(\mathcal{I}_v > 2\sigma) \leq \mathbb{P}_f(S_v > \sigma) + \mathbb{P}_f(|S_v - \mathcal{I}_v| > \sigma).$$

Using (4.5.15), (4.5.16) and applying Markov inequality, we get

$$(4.5.17) \quad \mathbb{P}_f(\mathcal{I}_v > 2\sigma) \leq T^{D_b} \sqrt{\Gamma_v} \sigma^{-1/2} e^{-g_h(v)} + C_h v^{D_b+1} T^{-1} \sigma^{-1}.$$

Note  $\mathbb{A} = \left\{ \int_{\hat{U}_n} Z_{h,\theta}(u) du < \frac{\delta^{D_b}}{2} \right\}$ . Since  $\mathcal{Q}_v \leq 1$  we obtain for any  $\delta > 0$  and  $\sigma > 0$

$$\begin{aligned} \mathbb{E}_f \mathcal{Q}_v &= \mathbb{E}_f [\mathcal{Q}_v \mathbb{I}_{\mathbb{A}} + \mathcal{Q}_v \mathbb{I}_{\mathcal{I}_v > 2\sigma, \mathbb{A}^c} + \mathcal{Q}_v \mathbb{I}_{\mathcal{I}_v \leq 2\sigma, \mathbb{A}^c}] \\ &\leq \mathbb{P}_f \left( \int_{\hat{U}_n} Z_{h,\theta}(u) du < \frac{\delta^{D_b}}{2} \right) + \mathbb{P}_f(\mathcal{I}_v > 2\sigma) + 4\delta^{-D_b} \sigma. \end{aligned}$$

Using to Lemma 8.3 and the inequality (4.5.17), we have

$$\mathbb{E}_f \mathcal{Q}_v \leq 2C_h \delta + T^{D_b} \sqrt{\Gamma_v} \sigma^{-1/2} e^{-g_h(v)} + C_h v^{D_b+1} T^{-1} \sigma^{-1} + 4\delta^{-D_b} \sigma.$$

Choosing  $T = \left\lceil \exp \left\{ \frac{1}{2D_b} g_h(v) \right\} \right\rceil + 1$ ,  $\sigma = \exp \left\{ -\frac{1}{3D_b} g_h(v) \right\}$  and  $\delta = \exp \left\{ -\frac{1}{6D_b^2} g_h(a) \right\}$ , we obtain

$$\mathbb{E}_f \mathcal{Q}_v \leq \left[ 2C_h + 2^{D_b} \sqrt{\Gamma_v} + C_h v^{D_b+1} + 4 \right] \exp \left\{ -\frac{1}{6D_b^2} g_h(a) \right\}.$$

Simplest algebra shows that  $\sqrt{\Gamma_v} \leq (2v+2)^{\frac{D_b}{2}}$ , we get

$$(4.5.18) \quad \mathbb{E}_f \mathcal{Q}_v \leq \left[ v^{D_b+1} + 2^{D_b} (2v+2)^{\frac{D_b}{2}} + 6 \right] C_h \exp \left\{ -\frac{1}{6D_b^2} g_h(a) \right\},$$

Note that if the event  $G$  is realized then  $\hat{U}_n(r) \subseteq U_n(r)$  and  $U_n(a) = \bigcup_{j=0}^{\infty} \Gamma_{a+j}$ . we obtain in view of (4.5.18)

$$\begin{aligned} \mathbb{E}_f \int_{\hat{U}_n \cap (\|u\|_1 > a)} \|u\|_1 z_h(u) \mathbb{I}_G du &\leq \sum_{j=0}^{\infty} (a+j+1) \mathbb{E}_f \mathcal{Q}_{a+j} \\ &= \Sigma^*(M(f)) a \mathcal{B}_a \mathcal{C}_h \exp \left\{ -\frac{1}{6D_b^2} g_h(a) \right\}. \end{aligned}$$

where we have put  $\mathcal{B}_a = a^{D_b+1} + 2^{D_b}(2a+2)^{\frac{D_b}{2}} + 6 + D_b \left( a^{D_b} + (2a+2)^{\frac{D_b}{2}-1} \right)$ .

*Third step: Conclusion.*

In view of Assertion 2, choosing  $r = \varepsilon/4$  we get

$$(4.5.19) \quad \mathbb{E}_f \int_{\hat{U}_n \cap (\|u\|_1 > \varepsilon/4)} \|u\|_1 z_h(u) \mathbb{I}_G du \leq \frac{\varepsilon}{4} \Sigma^*(M(f)) \mathcal{B}_{\varepsilon/4} \mathcal{C}_h e^{-\frac{1}{6D_b^2} g_h(\varepsilon/4)}.$$

Using the Markov inequality, we have in view of (4.5.19)

$$\mathbb{P}_f \left( \int_{\hat{U}_n \cap (\|u\|_1 > \varepsilon/4)} \|u\|_1 z_h(u) du > \frac{\varepsilon}{8}, G \right) \leq 2 \Sigma^*(M(f)) \mathcal{B}_{\varepsilon/4} \mathcal{C}_h e^{-\frac{1}{6D_b^2} g_h(\varepsilon/4)}.$$

The assertion of Pro follows now from the last inequality, Assertion 1 and the definitions of  $\mathcal{C}_h, g_h(\cdot)$  and the function  $\mathfrak{B}(\cdot, \cdot)$ . ■

### 4.5.3 Proof of Proposition 5

To prove the proposition it suffice to integrate the inequality obtained in Proposition 4 and to use the following lemma which will be extensively exploited in the sequel.

**Lemma 9.** *There exists  $\lambda > 0$  such that  $\forall n > 1$  and  $\forall h \in \mathcal{H}_n$ , we have*

$$\lambda_n(h) \geq \lambda.$$

where  $\lambda_n(h)$  is the smallest eigenvalue of the matrix

$$\mathcal{M}_{nh}(y) = \frac{1}{nh^d} \sum_{i=1}^n K^\top \left( \frac{X_i - y}{h} \right) K \left( \frac{X_i - y}{h} \right) \mathbb{I}_{V_h(y)}(X_i),$$

and  $K(z)$  is the  $D_b$ -dimensional vector of the monomials  $z^p$ ,  $p \in \mathcal{P}_b$ .



**Proof of Proposition 5.** In order to simplify the proof, let us introduce the following constants

$$(4.5.20) \quad c_1 = \frac{(1 + 6D_b^2)}{6A(f)D_b^2}, \quad c_2 = \frac{\lambda}{432M(f)D_b^3}.$$

By definition of  $A(f)$ ,  $M(f)$ ,  $\mathcal{B}(\cdot, \cdot)$  respectively in (4.1.11), (4.5.5) and  $A$ ,  $M$ , we have the following inequality  $\mathfrak{B}(A(f), M(f)) \leq \mathfrak{B}(A, M)$ . By integration of Proposition 4 and using Lemma 9, we get for any  $q \geq 1$  and  $f \in \mathbb{H}_d(\beta, L, M, A)$

$$(4.5.21) \quad \begin{aligned} & \mathbb{E}_f |\hat{f}^h(y) - f(y)|^q \mathbb{I}_G \\ &= \int_0^{+\infty} \eta^{q-1} \mathbb{P}_f \left( |\hat{f}^h(y) - f(y)| \geq \eta, G \right) d\eta \\ &= (nh^d)^{-q} \int_0^{+\infty} \eta^{q-1} \mathbb{P}_f \left( |\hat{f}^h(y) - f(y)| \geq \frac{\eta}{nh^d}, G \right) d\eta \\ &= (nh^d)^{-q} \left[ \int_0^{\frac{c_1}{c_2} \mathcal{N}_h} \eta^{q-1} d\eta \right. \\ & \quad \left. + \int_{\frac{c_1}{c_2} \mathcal{N}_h}^{+\infty} \eta^{q-1} \mathbb{P}_f \left( |\hat{f}^h(y) - f(y)| \geq \frac{\eta}{nh^d}, G \right) d\eta \right] \\ &\leq (nh^d)^{-q} \left[ \frac{c_1^q}{q c_2^q} \mathcal{N}_h^q + \int_{\frac{c_1}{c_2} \mathcal{N}_h}^{+\infty} \eta^{q-1} \mathfrak{B}(A, M) \mathcal{E}(h) e^{-c_2 \eta} d\eta \right]. \end{aligned}$$

In seeking this type of result (Control of the risk), the new point is to control the term  $\mathcal{E}(h)$  which can be very large compared to the rate  $(nh^d)^{-1}$ . It is a consequence that the bayesian estimator is non-linear. The bias-variance ratio  $\mathcal{N}_h$  (defined in (4.5.4)), which appears in  $\mathcal{E}(h)$ , allows to obtain similar results with linear estimators.

By calculation of the integral in (4.5.21),  $b_h$  and  $\mathcal{N}_h$  respectively defined in (4.5.3) and (4.5.4), the assertion of Proposition 5 is proved. ■

#### 4.5.4 Proof of Theorem 11

By definition of  $\bar{h} = (Ln)^{-\frac{1}{\beta+d}}$ , we have

$$(4.5.22) \quad Ldn\bar{h}^{\beta+d} = d, \quad (n\bar{h}^d)^{-q} = L^{\frac{qd}{\beta+d}} \varphi_n^q(\beta).$$

We take  $\hat{M} = M$  and  $\hat{A} = A$  (known in minimax case), by definition of  $\bar{f}^h$  in (4.1.7), (4.1.8), therefore Proposition 5 can be obtained by replacing  $\hat{f}$  by  $\bar{f}$ . Indeed the event  $G$ , defined in Section 4.5.1, is not necessary when the set of coefficients  $\Theta(A, M)$ , in the definition of

$\bar{f}$  (see (4.1.6)), is not random. In this case, the reasoning of proofs of Propositions 4 and 5 remains true for the estimator  $\bar{f}$ , then Proposition 5 can be applied for the estimator  $\bar{f}$  with same constants. Using (4.5.22), we obtain

$$\mathbb{E}_f |\bar{f}^h(y) - f(y)|^q \leq C_q^* d L^{\frac{qd}{\beta+d}} \varphi_n^q(\beta).$$

where  $\lambda$  is defined in the Lemma 5 in Section 4.5.3. The theorem 11 is proved. ■

### 4.5.5 Proof of Theorem 14

This Proof is based on the Lepski scheme developed by Lepski [1991] and adapted for the bandwidth selection by Lepski, Mammen, and Spokoiny [1997]. We start the proof with formulating auxiliary Lemmas whose proofs are given in Appendix (Section 4.7). Define  $J_1 = Ld(1 + 6D_b^2)/6A(f)D_b^2$  and

$$h^* = \left[ \frac{c(1 + (b - \beta) \ln n)}{n} \right]^{\frac{1}{\beta+d}}, \quad c < \frac{1 \wedge qdJ_1^{-1}}{(b + d)(\beta + d)}.$$

and let the integer  $\kappa$  be defined as follows.

$$(4.5.23) \quad 2^{-\kappa} h_{\max} \leq h^* < 2^{-\kappa+1} h_{\max}.$$

The definitions of  $h^*$  and  $\kappa$  (4.5.23) imply the following Lemmas.

**Lemma 10.**

$$\mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^q \mathbb{I}_G \leq \bar{C}_q \frac{(1 + k \ln 2)^q}{(nh_k^d)^q}, \quad \forall k \geq \kappa,$$

$$\text{where } \bar{C}_q = C_q^* \frac{(\beta + d)(1 \wedge qdJ_1^{-1})}{(\beta + d - 1)(Ld)^{-1}}.$$

**Lemma 11.** For any  $f \in \mathbb{H}_d(\beta, L, M, A)$  and any  $k \geq \kappa + 1$

$$\mathbb{P}_f(\hat{k} = k, G) \leq J_2 \mathfrak{B}(A, M) \exp \{ J_1 n (h^*)^{\beta+d} \} 2^{-(k-1)(8qd+4)},$$

where  $J_2 = (1 - 2^{-(8qd+4)})^{-1}$ .

**Lemma 12.** There exists a universal constant  $\vartheta > 0$  such that

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, M, A)} \exp \left\{ \frac{An^{\frac{b}{b+d}}}{16M\vartheta^2 D_b^2} \right\} \mathbb{P}_f(G^c) = 0.$$

**Proof of Theorem 14.** We decompose the risk as follows

$$\begin{aligned}
 & \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_G \\
 & \leq \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_{\hat{k} \leq \kappa, G} + \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_{\hat{k} > \kappa, G} \\
 (4.5.24) \quad & = R_1(f) + R_2(f).
 \end{aligned}$$

First we control  $R_1$ . Obviously

$$|\hat{f}^{(\hat{k})}(y) - f(y)| \leq |\hat{f}^{(\hat{k})}(y) - \hat{f}^{(\kappa)}(y)| + |\hat{f}^{(\kappa)}(y) - f(y)|.$$

Note that the realization of the event  $G$  implies  $\hat{M} \leq 3M(f)/2$ . This together with the definition of  $\hat{k}$  yields

$$|\hat{f}^{(\hat{k})}(y) - \hat{f}^{(\kappa)}(y)| \mathbb{I}_{\hat{k} \leq \kappa, G} \leq C s_n(\kappa), \quad s_n(k) = (1 + k \ln 2)^q (n h_k^d)^{-q},$$

where  $C = 288 M D_b^3 \lambda^{-1} (32 q d + 16)$ . In view of Lemma 10 we also get

$$\mathbb{E}_f |\hat{f}^{(\kappa)}(y) - f(y)|^q \leq \bar{C}_q s_n(\kappa).$$

Noting that the right hand side of the obtain inequality is independent of  $f$  and taking into account the definition of  $\kappa$  and  $h^*$  we obtain

$$(4.5.25) \quad \limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-q}(\beta) R_1(f) < \infty.$$

Now let us bounded from above  $R_2$ . Applying Cauchy-Schwartz inequality we have in view of Lemma 11

$$\begin{aligned}
 R_2(f) &= \sum_{k > \kappa}^{k_n} \mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^q \mathbb{I}_{[\hat{k}=k, G]} \\
 &\leq \sum_{k > \kappa} (\mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^{2q})^{1/2} \sqrt{\mathbb{P}_f \{\hat{k} = k, G\}} \\
 (4.5.26) \quad &= \Delta(h^*) \sum_{k > \kappa} (\mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^{2q})^{1/2} 2^{-(k-1)(4qd+2)},
 \end{aligned}$$

where we have put  $\Delta(h^*) = J_2 \mathfrak{B}(A, M) \exp \{J_1 n (h^*)^{\beta+d}\}$ . We obtain from Lemma 10 and (4.5.26)

$$(4.5.27) \quad R_2(f) \leq J_3 (n h_{\max}^d)^{-q} \exp \{J_1 n (h^*)^{\beta+d}\},$$

where

$$J_3 = J_2 \mathfrak{B}(A, M) 2^{4qd+2} \bar{C}_{2q}^{1/2} \sum_{s \geq 0} (1 + s \ln 2)^q 2^{-3sdq-2}.$$

It remains to note that the definition of  $h^*$  implies that

$$\limsup_{n \rightarrow \infty} \phi_n^{-q}(\beta) (nh_{\max}^d)^{-q} \exp \{J_1 n (h^*)^{\beta+d}\} < \infty$$

and that the right hand side of (4.5.27) is independent of  $f$ . Thus, we have

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-q}(\beta) R_2(f) < \infty.$$

that yields together with (4.5.24) and (4.5.25)

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-q}(\beta) \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_G < \infty.$$

To get the assertion of the theorem it suffices to show that

$$(4.5.28) \quad \limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-q}(\beta) \mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_{G^c} < \infty.$$

Note that  $\hat{f}^{(\hat{k})}(y) \leq 4\hat{M}$  in view of (4.3.6). Note also that the local least square estimator  $\tilde{\delta}$  is linear function of observation  $Y^{(n)}$  and, moreover  $0 \leq Y_i \leq M, i = 1, \dots, n$ . This together with the definition of  $\hat{M}$ , (expression (4.3.5)) allows us to state that there exist  $0 < J_4 < +\infty$  such that  $|\hat{f}^{(\hat{k})}(y) - f(y)| \leq J_4 M$ . Here we also have taken into account that  $\|f\|_\infty \leq M$ .

Finally we obtain

$$\mathbb{E}_f |\hat{f}^{(\hat{k})}(y) - f(y)|^q \mathbb{I}_{G^c} \leq J_4^q M^q \mathbb{P}_f \{G^c\}.$$

and (4.5.28) follows now from Lemma 11. ■

## 4.6 Proofs of lower bounds

The proofs of Theorems 10 and 13 are based on the following proposition.

Put  $\phi_n(\gamma) = [n^{-1}(1 + (b - \gamma) \ln n)]^{\frac{\gamma}{\gamma+d}}$ ,  $\gamma \in (0, b]$  and let

$$\begin{aligned} R_n^{(q)}(\tilde{f}, v) &= \sup_{f \in \mathbb{H}_d(\alpha, L, M, A)} \mathbb{E}_f \left[ \phi_n^{-q}(\alpha) |\tilde{f}(y) - f(y)|^q \right] \\ &\quad + \sup_{f \in \mathbb{H}_d(\beta, L, M, A)} \mathbb{E}_f \left[ n^{-vq} \phi_n^{-q}(\beta) |\tilde{f}(y) - f(y)|^q \right]. \end{aligned}$$

where  $v \geq 0$  and  $\alpha, \beta \in (0, b]^2$ .

**Proposition 6.** *Let  $\Psi$  be admissible family of normalizations such that*

$$\psi_n(\alpha) / \phi_n(\alpha) \xrightarrow{n \rightarrow \infty} 0.$$

*Then, for any  $0 \leq v < (\beta - \alpha) / (\beta + 1)(\alpha + 1)$*

$$\liminf_{n \rightarrow \infty} \inf_{\tilde{f}} R_n^{(q)}(\tilde{f}, v) > 0.$$

The proof is given in section 4.6.3.

### 4.6.1 Proof of Theorem 10

Using the proposition 6 for  $\beta = \alpha$ , we have to choose  $v = 0$  and one gets

$$\begin{aligned} R_{n,q}[\mathbb{H}_d(\beta, L, M, A)] &= R_n^{(q)}(\tilde{f}, 0) \\ &= \sup_{f \in \mathbb{H}_d(\alpha, L, M, A)} \mathbb{E}_f \left[ n^{-q \frac{\alpha}{\alpha+1}} |\tilde{f}(y) - f(y)|^q \right] > 0, \quad \forall \tilde{f}. \end{aligned}$$

■

### 4.6.2 Proof of Theorem 13

**I.** To proof of the first assertion of the theorem it suffices to consider the family  $\{v_n(\beta)\}_{\beta \in (0, b]}$ , where  $v_n(\alpha) = \varphi_n(\alpha)$  and  $v_n(\beta) = 1$  for any  $\beta \neq \alpha$ . The corresponding attainable estimator is the estimator being minimax on  $\mathbb{H}_d(\alpha, L, M, A)$ .

**II.** Let us consider the family  $\{\phi_n(\beta)\}_{\beta \in [0, b]}$ , which is admissible in view of Theorem 14. First, we note that  $\gamma = b$  is not possible since  $\phi_n(b) = \varphi_n(b)$  the minimax rate of convergence on  $\mathbb{H}_d(b, L, M, A)$ .

Thus we assume that  $\gamma$  satisfying (4.3.2) belongs to  $\delta \in ]0, b[$ . Let  $\hat{f}^\Psi$  be a  $\Psi^{(n)}$ -attainable estimator. Since  $\psi_n(\alpha)/\phi_n(\alpha) \rightarrow 0, n \rightarrow \infty$  in view of (4.3.2) then obviously

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\gamma, L, M, A)} \mathbb{E}_f \left[ \phi_n^{-q}(\gamma) |\hat{f}^\Psi(y) - f(y)|^q \right] = 0.$$

Therefore, applying Proposition 6 with  $v = 0$  we have for any  $\beta < \gamma$

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, M, A)} \mathbb{E}_f \left[ \phi_n^{-q}(\beta) |\hat{f}^\Psi(y) - f(y)|^q \right] > 0.$$

We conclude that necessarily  $\psi_n(\beta) \gtrsim \phi_n(\beta)$  for any  $\beta < \gamma$ .

Moreover for any  $\beta > \gamma$  applying Proposition 6 with an arbitrary  $0 \leq v < (\beta - \gamma)/(\beta + 1)(\gamma + 1)$  we obtain that

$$\psi_n(\beta) \gtrsim n^v \phi_n(\beta), \quad \beta > \gamma.$$

It remains to note that the form of rate of convergence proved in Theorem 10 implies that

$$\phi_n(\gamma)/\psi_n(\gamma) = o\left([\ln n]^{\frac{\gamma}{\gamma+d}}\right).$$

■

### 4.6.3 Proof of Proposition 6

Let  $\varkappa > 0$  the parameter whose choice will be done later. Put

$$h = \left( \varkappa \frac{1 + (\beta - \alpha) \ln n}{n} \right)^{\frac{1}{\alpha + d}}.$$

Later on without loss of generality we will assume that  $L > 1$ .

Consider the functions:  $f_0 \equiv 1$  and

$$f_1(x) = 1 - (L - 1) \varkappa^{\frac{\alpha}{\alpha + d}} \phi_n(\alpha) F \left( \frac{x_1 - y_1}{h}, \dots, \frac{x_d - y_d}{h} \right), \quad x \in [0, 1]^d.$$

Here  $F$  is a compactly supported positive function belonging to  $\mathbb{H}_d(\alpha, 1, M, A)$  such that  $F(0) = 1 = \max_x F(x)$ .

It is easily seen that  $f_1 \in \mathbb{H}_d(\alpha, L, M, A)$ . Therefore, we have

$$\begin{aligned} R_n^{(q)}(\tilde{f}, v) &\geq \mathbb{E}_0 \left| n^{-v} \phi_n^{-1}(\beta) (\tilde{f}(y) - 1) \right|^q + \mathbb{E}_1 \left| \phi_n^{-1}(\alpha) (\tilde{f}(y) - f_1(y)) \right|^q \\ &\geq \mathbb{E}_0 \left| n^{-v} \phi_n^{-1}(\beta) (\tilde{f}(y) - 1) \right|^q + \mathbb{E}_1 \left| \phi_n^{-1}(\alpha) (\tilde{f}(y) - 1) + z \right|^q, \end{aligned}$$

where  $z = (L - 1) \varkappa^{\frac{\alpha}{\alpha + 1}} F(0)$ . Set

$$\tilde{\lambda} = \phi_n^{-1}(\alpha) (1 - \tilde{f}(y)), \quad \varsigma_n = n^{-v} \frac{\phi_n(\alpha)}{\phi_n(\beta)} = n^{-v} \left( \frac{\ln n}{n} \right)^{-\varrho},$$

where  $\varrho = \frac{\beta - \alpha}{(\beta + 1)(\alpha + 1)}$ . We get

$$\begin{aligned} R_n^{(q)}(\tilde{f}, v) &\geq \mathbb{E}_0 |\varsigma_n \tilde{\lambda}|^q + \mathbb{E}_1 |z - \tilde{\lambda}|^q \\ &\geq \mathbb{E}_0 |\varsigma_n \tilde{\lambda}|^q \mathbb{I}_{\{|\tilde{\lambda}| > z/2\}} + \mathbb{E}_1 |z - \tilde{\lambda}|^q \mathbb{I}_{\{|\tilde{\lambda}| \leq z/2\}} \\ (4.6.1) \quad &\geq \mathbb{E}_0 |\varsigma_n \frac{z}{2}|^q \mathbb{I}_{\{|\tilde{\lambda}| > z/2\}} + \mathbb{E}_1 |\frac{z}{2}|^q \mathbb{I}_{\{|\tilde{\lambda}| \leq z/2\}}. \end{aligned}$$

Noting that  $f_1 \leq f_0$ , since  $F$  is positive, and putting  $c_n(Y^{(n)}) = \mathbb{I}_{\{|\tilde{\lambda}| > z/2\}}$  we obtain

$$\begin{aligned} R_n^{(q)}(\tilde{f}, v) &\geq \varsigma_n^q \frac{z^q}{2^q} \frac{\prod_{i=1}^n f_1(X_i)}{\prod_{i=1}^n f_1(X_i)} \int_0^{f_1(X_1)} \dots \int_0^{f_1(X_n)} c_n(x) dx_1 \dots dx_n \\ (4.6.2) \quad &+ \frac{z^q}{2^q} \frac{1}{\prod_{i=1}^n f_1(X_i)} \int_0^{f_1(X_1)} \dots \int_0^{f_1(X_n)} 1 - c_n(x) dx_1 \dots dx_n. \end{aligned}$$

We have

$$\begin{aligned} \prod_{i=1}^n f_1(X_i) &= \prod_{i=1}^n \left( 1 - (L - 1) \varkappa^{\frac{\alpha}{\alpha + d}} \phi_n(\alpha) F \left( \frac{X_i - y}{h} \right) \right) \\ (4.6.3) \quad &\geq \left( 1 - (L - 1) \varkappa^{\frac{\alpha}{\alpha + d}} \phi_n(\alpha) \right)^{nh^d} \geq e^{-(L-1)\varkappa} n^{-(L-1)\varkappa(\beta-\alpha)}. \end{aligned}$$

We obtain in view of (4.6.2) and (4.6.3)

$$\begin{aligned}
R_n^{(q)}(\tilde{f}, v) &\geq \varsigma_n^q \frac{z^q}{2^q} e^{-(L-1)\varkappa} n^{-(L-1)\varkappa(\beta-\alpha)} \\
&\quad \times \frac{1}{\prod_{i=1}^n f_1(X_i)} \int_0^{f_1(X_1)} \dots \int_0^{f_1(X_n)} c_n(x) dx_1 \dots dx_n \\
&\quad + \frac{z^q}{2^q} \frac{1}{\prod_{i=1}^n f_1(X_i)} \int_0^{f_1(X_1)} \dots \int_0^{f_1(X_n)} 1 - c_n(x) dx_1 \dots dx_n \\
&\geq \frac{z^q}{2^q} \left(1 \wedge \varsigma_n^q e^{-(L-1)\varkappa} n^{-(L-1)\varkappa(\beta-\alpha)}\right).
\end{aligned}$$

Case 1:  $\beta = \alpha$ . Choosing  $\varkappa = 1$ , and noting that  $\varsigma_n = 1$  and  $\prod_{i=1}^n f_1(X_i) \geq e^{-(L-1)}$ , we deduce from (4.6.2) that yields:

$$\inf_{\tilde{f}} R_n^{(q)}(\tilde{f}, v) \geq \frac{(L-1)^q}{2^q} e^{-(L-1)} > 0.$$

Case 2:  $\beta > \alpha$ . Put

$$\varkappa = \frac{q(\varrho - v) - t_n}{1 + (L-1)(\beta - \alpha)} > 0, \quad t_n = \frac{q}{\ln n} \ln \frac{1}{(1 + (\beta - \alpha) \ln n)^{-\varrho}} \xrightarrow{n \rightarrow \infty} 0.$$

This choice provides us with the following bound

$$\begin{aligned}
\varsigma_n^q e^{-(L-1)\varkappa} n^{-(L-1)\varkappa(\beta-\alpha)} &= (1 + (\beta - \alpha) \ln n)^{-q\varrho} e^{-(L-1)\varkappa} n^{q(\varrho-v) - (L-1)\varkappa(\beta-\alpha)} \\
&\geq (1 + (\beta - \alpha) \ln n)^{-q\varrho} e^{-\frac{2}{3}q(L-1)} n^{t_n} \geq e^{-\frac{2}{3}q(L-1)}.
\end{aligned}$$

This yields

$$\inf_{\tilde{f}} R_n^{(q)}(\tilde{f}, v) \geq \frac{(L-1)^q \varkappa^{\frac{q}{\alpha+1}}}{2^q} e^{-\frac{2}{3}q(L-1)} > 0.$$

■

## 4.7 Appendix

**Proof of Lemma 8** Later on without loss generality we will suppose that  $nh^d \in \mathbb{N}^*$ . In order to simplify understanding of this proof, we note the approximation polynomial  $\mathcal{A}_u^i = f_{\theta+u(nh^d)^{-1}}(X_i)$ ,  $i = 1, \dots, n$  for all  $u \in U_n$ .

1. Note that for  $u \in U_n$

$$(4.7.1) \quad \mathbb{E}_f Z_{h,\theta}(u) \leq \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{f(X_i)} \leq e^{\mathcal{N}_h/A(f)}.$$

This result is the consequence of the definition of  $Z_{h,\theta}$  in (4.5.9) and the following calculation

$$\mathbb{E}_f \mathbb{I}_{[Y_i \leq \mathcal{A}_u^i]} = \mathbb{P}_f(Y_i \leq \mathcal{A}_u^i) = 1 \wedge \frac{\mathcal{A}_u^i}{f(X_i)}.$$

In (4.7.1), the second inequality is obtained with classical inequality  $1 + \rho \leq e^\rho$ ,  $\rho \in \mathbb{R}$  and recall that  $f_\theta(x) \geq f(x)$ .

$$\prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{f(X_i)} = \prod_{i: X_i \in V_h(y)} \left( 1 + \frac{\mathcal{A}_0^i - f(X_i)}{f(X_i)} \right) \leq \exp \{b_h \times nh^d/A(f)\}$$

Case 1: If  $\|u_1 - u_2\|_1 \geq 1$ , the inequality (4.7.1) allows to get

$$\mathbb{E}_f |Z_{h,\theta}(u_1) - Z_{h,\theta}(u_2)| \leq \mathbb{E}_f Z_{h,\theta}(u_1) + \mathbb{E}_f Z_{h,\theta}(u_2) \leq 2e^{\mathcal{N}_h/A(f)} \|u_1 - u_2\|_1.$$

Case 2: Assume now that  $\|u_1 - u_2\|_1 < 1$  and introduce the random events

$$\begin{aligned} F_1 &= \{ \forall i = 1, \dots, n : Y_i \leq \mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i \}, \\ F_2 &= \{ \forall i = 1, \dots, n : Y_i \leq \mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i \} \\ &\quad \cap \{ \exists i : Y_i > \mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i \}, \\ F_3 &= \{ \exists i : Y_i > \mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i \}. \end{aligned}$$

We have used the following notations:  $a \wedge b = \min(a, b)$  and  $a \vee b = \max(a, b)$ ,  $a, b \in \mathbb{R}$ . For any  $(u_1, u_2) \in U_n^2$ , we have

$$\begin{aligned} \mathbb{E}_f |Z_{h,\theta}(u_1) - Z_{h,\theta}(u_2)| &= \mathbb{E}_f |Z_{h,\theta}(u_1) - Z_{h,\theta}(u_2)| \mathbb{I}_{[F_1]} \\ &\quad + \mathbb{E}_f |Z_{h,\theta}(u_1) - Z_{h,\theta}(u_2)| \mathbb{I}_{[F_2]} + \mathbb{E}_f |Z_{h,\theta}(u_1) - Z_{h,\theta}(u_2)| \mathbb{I}_{[F_3]} \\ &= \mathcal{K}_1 + \mathcal{K}_2 + \mathcal{K}_3. \end{aligned}$$

The following bound will be extensively exploited in the sequel.

$$f_v(x) \geq 2v_{0,\dots,0} - \|v\|_1 \geq 0.25A(f), \quad \forall v \in \Theta(A(f)/4, 9M(f)), \quad x \in [0, 1]^d.$$



**Control of  $\mathcal{K}_1$ .**

$$(4.7.2) \quad \mathcal{K}_1 = \left| \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{\mathcal{A}_{u_1}^i} - \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{\mathcal{A}_{u_2}^i} \right| \mathbb{P}_f\{F_1\},$$

and

$$(4.7.3) \quad \mathbb{P}_f\{F_1\} = \prod_{i: X_i \in V_h(y)} \mathbb{P}_f\{Y_i \leq \mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i\} \leq \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i}{f(X_i)}.$$

Therefore,

$$(4.7.4) \quad \begin{aligned} \mathcal{K}_1 &\leq \left( 1 - \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i}{\mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i} \right) \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{f(X_i)} \\ &\leq e^{\mathcal{N}_h/A(f)} \left( 1 - \exp \left\{ \sum_{i: X_i \in V_h(y)} \ln \frac{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i}{\mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i} \right\} \right). \end{aligned}$$

Remember that  $|\mathcal{A}_{u_1}^i - \mathcal{A}_{u_2}^i| \leq (nh^d)^{-1} \|u_1 - u_2\|_1$  and  $\mathcal{A}_u^i \geq A(f)/4$ . Let us give the following calculation with inequality of finite increments

$$\ln \frac{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i}{\mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i} = - |\ln \mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i - \ln \mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i| \geq - \frac{(nh^d)^{-1} \|u_1 - u_2\|_1}{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i}$$

Using last inequalities, (4.7.2), (4.7.3), (4.7.4) and the well known inequality  $1 - e^{-\rho} \leq \rho$ , we have

$$\mathcal{K}_1 \leq \frac{1}{A(f)} e^{\mathcal{N}_h/A(f)} \|u_1 - u_2\|_1.$$

**Control of  $\mathcal{K}_2$ .** Put

$$\begin{aligned} F_2 &= \{ \forall i = 1, \dots, n : Y_i \leq \mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i \} \\ &\quad \setminus \{ \forall i = 1, \dots, n : Y_i \leq \mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i \} \\ &= G \setminus F_1. \end{aligned}$$

and define

$$\begin{aligned} \mathcal{G}_1 &= \{ X_i \in V_h(y) : \mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i < f(X_i) \}, \\ \mathcal{G}_2 &= \{ X_i \in V_h(y) : \mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i < f(X_i) \}. \end{aligned}$$

Note that  $F_1 \subseteq G$  and, therefore,

$$\begin{aligned} \mathcal{K}_2 &\leq \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i} (\mathbb{P}_f\{G\} - \mathbb{P}_f\{F_1\}) \\ &= \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i} \left( \prod_{i: X_i \in \mathcal{G}_1} \frac{\mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i}{f(X_i)} - \prod_{i: X_i \in \mathcal{G}_2} \frac{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i}{f(X_i)} \right) \end{aligned}$$

The definition of  $\mathcal{G}_2$  implies

$$\prod_{i: X_i \in V_h(y)} \frac{1}{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i} \leq \prod_{i: X_i \in \mathcal{G}_2} \frac{1}{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i} \prod_{i: X_i \in \mathcal{G}_2^c} \frac{1}{f(X_i)}$$

Since  $\mathcal{G}_1 \subseteq \mathcal{G}_2$ ,  $\|u_1 - u_2\|_1 < 1$  and  $|f_u(x)| \leq \|u\|_1$ ,  $\forall x \in [0, 1]^d$ ,  $\forall u \in U_n$ , using the last inequality, we obtain

$$\begin{aligned} \mathcal{K}_2 &\leq \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{f(X_i)} \left( \prod_{i: X_i \in \mathcal{G}_2} \frac{\mathcal{A}_{u_1}^i \vee \mathcal{A}_{u_2}^i}{\mathcal{A}_{u_1}^i \wedge \mathcal{A}_{u_2}^i} - 1 \right) \\ &\leq 4D_b e^{1+\mathcal{N}_h/A(f)} \|u_1 - u_2\|_1 / A(f). \end{aligned}$$

It remains to note that  $\mathcal{K}_3 = 0$  and the first assertion of the lemma is proved with the bounds of  $\mathcal{K}_1$  and  $\mathcal{K}_2$ .

**2.** For any  $u \in U_n$ , since the random variables  $(Y_i)_i$  are independent we have,

$$\mathbb{E}_f Z_{n,\theta}^{1/2}(u) = \prod_{i: X_i \in V_h(y)} \sqrt{\frac{\mathcal{A}_0^i}{\mathcal{A}_u^i}} \mathbb{P}_f \{Y_i \leq \mathcal{A}_u^i\}.$$

For any  $i$ , we have

$$\sqrt{\frac{\mathcal{A}_0^i}{\mathcal{A}_u^i}} \mathbb{P}_f \{Y_i \leq \mathcal{A}_u^i\} = \sqrt{\frac{\mathcal{A}_0^i}{\mathcal{A}_u^i}} \left[ 1 \wedge \frac{\mathcal{A}_u^i}{f(X_i)} \right] \leq \frac{\mathcal{A}_0^i}{f(X_i)} \left[ \frac{f(X_i)}{\sqrt{\mathcal{A}_0^i} \sqrt{\mathcal{A}_u^i}} \wedge \frac{\sqrt{\mathcal{A}_u^i}}{\sqrt{\mathcal{A}_0^i}} \right].$$

Remind that in view of (4.5.8)  $f_\theta(x) \geq f(x)$  and  $0 < f_\theta(x) \leq 3M(f)$  for  $x \in V_h(y)$ .

Moreover, for  $u \in U_n = nh^d(\Theta(A(f)/4, 9M(f)) - \theta)$ ,  $0 < f_{\theta+u(nh^d)-1}(x) \leq 9M(f)$ . Thus for all  $i : X_i \in V_h(y)$ ,

$$\sqrt{\frac{\mathcal{A}_0^i}{\mathcal{A}_u^i}} \mathbb{P}_f \{Y_i \leq \mathcal{A}_u^i\} \leq \frac{\mathcal{A}_0^i}{f(X_i)} \left[ \frac{\sqrt{\mathcal{A}_0^i}}{\sqrt{\mathcal{A}_u^i}} \wedge \frac{\sqrt{\mathcal{A}_u^i}}{\sqrt{\mathcal{A}_0^i}} \right] \leq \frac{\mathcal{A}_0^i}{f(X_i)} \left[ 1 - \frac{|\mathcal{A}_0^i - \mathcal{A}_u^i|}{9M(f)} \right]^{1/2}.$$

The last inequality implies

$$\begin{aligned}
 \mathbb{E}_f Z_{n,\theta}^{1/2}(u) &\leq \prod_{i: X_i \in V_h(y)} \frac{\mathcal{A}_0^i}{f(X_i)} \sqrt{1 - \frac{|f_{u(nh^d)^{-1}}(X_i)|}{9M(f)}} \\
 (4.7.5) \quad &\leq e^{\mathcal{N}_h/A(f)} \exp \left\{ -\frac{1}{18M(f) nh^d} \sum_{i: X_i \in V_h(y)} |f_u(X_i)| \right\}.
 \end{aligned}$$

It remains to show

$$(4.7.6) \quad \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} |f_u(X_i)| \geq \lambda_n(h) D_b^{-1} \|u\|_1.$$

Let us remember that  $u = (u_p, p \in \mathcal{P}_b)$  (where  $\mathcal{P}_b$  is defined in (4.1.5)). First, we get from the definition of  $f_u$

$$f_u(x) = u K^\top \left( \frac{x-y}{h} \right) = K \left( \frac{x-y}{h} \right) u^\top, \quad \forall x \in [0, 1]^d,$$

and, therefore,

$$\frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} |f_u(X_i)| = \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \left| u K^\top \left( \frac{X_i - y}{h} \right) \right|.$$

Assume  $u \neq 0$  and put  $v = u/\|u\|_1$ . Noting that  $|f_v(x)| \leq 1$ ,  $\forall x \in [0, 1]^d$ , we have

$$\begin{aligned}
 &\frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} |f_u(X_i)| \\
 &\geq \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \left| u K^\top \left( \frac{X_i - y}{h} \right) \right| |f_v(X_i)| \\
 &= \frac{1}{\|u\|_1 nh^d} \sum_{i: X_i \in V_h(y)} \left| u K^\top \left( \frac{X_i - y}{h} \right) K \left( \frac{X_i - y}{h} \right) u^\top \right| \\
 &\geq \frac{1}{\|u\|_1} \left| u \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} K^\top \left( \frac{X_i - y}{h} \right) K \left( \frac{X_i - y}{h} \right) u^\top \right|.
 \end{aligned}$$

The bound (4.7.6) follows now from Lemma 9. The assertion of the lemma follows from (4.7.5) and (4.7.6).

**3.** In view of Lemma 8.1, we have

$$(4.7.7) \quad \mathbb{E}_f |Z_{h,\theta}(u) - Z_{h,\theta}(0)| \leq \mathcal{C}_h \|u\|_1, \quad u \in U_n.$$

Taking into account that  $Z_{h,\theta}(0) = 1$  we obtain applying (4.7.7), Fubini's theorem and Markov inequality

$$\begin{aligned}
& \mathbb{P}_f \left\{ \int_0^\delta \cdots \int_0^\delta Z_{h,\theta}(v) dv < \frac{1}{2} \delta^{D_b} \right\} \\
&= \mathbb{P}_f \left\{ \int_0^\delta \cdots \int_0^\delta (Z_{h,\theta}(v) - Z_{h,\theta}(0)) dv < -\frac{1}{2} \delta^{D_b} \right\} \\
&\leq \mathbb{P}_f \left\{ \int_0^\delta \cdots \int_0^\delta |Z_{h,\theta}(v) - Z_{h,\theta}(0)| dv > \frac{1}{2} \delta^{D_b} \right\} \\
&\leq 2\delta^{-D_b} \int_0^\delta \cdots \int_0^\delta \mathbb{E}_f |Z_{h,\theta}(v) - Z_{h,\theta}(0)| dv \\
&\leq 2\mathcal{C}_h \delta
\end{aligned}$$

■

**Proof of Lemma 9** *First step:  $\mathcal{M}_{nh}(y)$  is a nonnegative positive matrix.*

Let  $\mathcal{H}_n, n > 1$  is defined in (4.5.1). First, we prove that

$$(4.7.8) \quad \inf_{h \in \mathcal{H}_n} \lambda_n(h) > 0, \quad \forall n > 1.$$

Suppose that  $\exists n_1 > 1, h_{n_1} \in \mathcal{H}_{n_1}$  such that  $\lambda_{n_1}(h_{n_1}) = 0$ . Recall that  $f_t(x) = tK(h^{-1}(x-y))$  for all  $t \in \mathbb{R}^{D_b}$  and note that  $\forall \tau \in \mathbb{R}^{D_b}$

$$\begin{aligned}
\tau^\top \mathcal{M}_{n_1 h_{n_1}}(y) \tau &= \frac{1}{nh_{n_1}^d} \sum_{i: X_i \in V_{h_{n_1}}(y)} \left[ \tau^\top K^\top \left( \frac{X_i - y}{h_{n_1}} \right) \right]^2 \\
&= \frac{1}{nh_{n_1}^d} \sum_{i: X_i \in V_{h_{n_1}}(y)} [f_\tau(X_i)]^2 \geq 0.
\end{aligned}$$

Since  $\lambda_{n_1}(h_{n_1})$  is the smallest eigenvalue of the matrix  $\mathcal{M}_{n_1 h_{n_1}}(y)$  the assumption  $\lambda_{n_1}(h_{n_1}) = 0$  implies that there exist  $\tau^*$  belonging to the unit sphere of  $\mathbb{R}^{D_b}$  such that

$$\frac{1}{nh_{n_1}^d} \sum_{i: X_i \in V_{h_{n_1}}(y)} [f_{\tau^*}(X_i)]^2 = 0.$$

It obviously implies that  $f_{\tau^*}(X_i) = 0$  for all  $X_i \in V_{h_{n_1}}(y)$ . It remains to note that  $nh_{n_1}^d \geq (b+1)^d$  since  $h_{n_1} \in \mathcal{H}_n$  and to apply the result obtained in Nemirovski [2000] (page 20). It yields  $\tau^* = 0$  and the obtained contradiction proves (4.7.8).

*Second step:*  $\mathcal{M}_{nh}(y) \xrightarrow{n \rightarrow \infty} \mathcal{M}$ .

Let  $\lambda_0$  be the smallest eigenvalue of the matrix

$$\mathcal{M} = \int_{[-1/2, 1/2]^d} K^\top(x) K(x) dx$$

whose general term is given by

$$\mathcal{M}_{p,q} = \prod_{j=1}^d \int_{-\frac{1}{2}}^{\frac{1}{2}} x_j^{p_j+q_j} dx_j, \quad 0 \leq |p|, |q| \leq b.$$

Let us prove that

$$(4.7.9) \quad \limsup_{n \rightarrow \infty} \sup_{h \in \mathcal{H}_n} |\lambda_n(h) - \lambda_0| = 0.$$

Put  $m = n^{1/d}$  and without loss of generality we will assume that  $m$  is integer. Remind that the general term of the matrix  $\mathcal{M}_{nh}(y)$  is given by

$$(\mathcal{M}_{nh}(y))_{p,q} = \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \prod_{j=1}^d \left( \frac{X_{i_j} - y_j}{h} \right)^{p_j+q_j}.$$

where  $X_{i_j} = i_j/m$  for all  $j = 1, \dots, d$  and  $X_i = (X_{i_1}, \dots, X_{i_d})$ . We get

$$\begin{aligned} & \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \prod_{j=1}^d \int_{i_j-1}^{i_j} \left( \frac{x_j/m - y_j}{h} \right)^{p_j+q_j} dx_j \\ & \leq \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \prod_{j=1}^d \left( \frac{X_{i_j} - y_j}{h} \right)^{p_j+q_j} \\ & \leq \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \prod_{j=1}^d \int_{i_j}^{i_j+1} \left( \frac{x_j/m - y_j}{h} \right)^{p_j+q_j} dx_j, \end{aligned}$$

It yields by change of variables that

$$\begin{aligned} \prod_{j=1}^d \int_{-\frac{1}{2}-2(nh^d)^{-1}}^{\frac{1}{2}} x_j^{p_j+q_j} dx_j & \leq \frac{1}{nh^d} \sum_{i: X_i \in V_h(y)} \prod_{j=1}^d \left( \frac{X_{i_j} - y_j}{h} \right)^{p_j+q_j} \\ (4.7.10) \quad & \leq \prod_{j=1}^d \int_{-\frac{1}{2}}^{\frac{1}{2}+2(nh^d)^{-1}} x_j^{p_j+q_j} dx_j, \end{aligned}$$

Note that  $nh^d \geq \ln^{\frac{1}{1+d}}(n)$  for any  $h \in \mathcal{H}_n$ . This together with (4.7.10) yields

$$\limsup_{n \rightarrow \infty} \sup_{h \in \mathcal{H}_n} \left| (\mathcal{M}_{nh}(y))_{p,q} - \mathcal{M}_{p,q} \right| = 0, \quad 0 \leq |p|, |q| \leq b.$$

The last result obviously imply (4.7.9).

*Third step: Conclusion.*

First we show that  $\lambda_0 > 0$ . Indeed,  $\forall \tau \in \mathbb{R}^{D_b}$

$$\tau^\top \mathcal{M} \tau = \int_{[-1/2, 1/2]^d} [f_\tau(x)]^2 dx \geq 0.$$

Since  $\lambda_0$  is smallest eigenvalue of the matrix  $\mathcal{M}$  the assumption  $\lambda_0 = 0$  would imply that there exists  $\tau^*$  belonging to the unit sphere of  $\mathbb{R}^{D_b}$  such that  $f_{\tau^*} \equiv 0$ . Since  $f_{\tau^*}$  is a polynomial the last identity is possible if and only if  $\tau^* = 0$ . The obtained contradiction shows that  $\lambda_0 > 0$ .

Next, note that in view of (4.7.9) there exists  $n_0$  such that  $\forall n > n_0$  and  $\forall h \in \mathcal{H}_n$ ,  $\lambda_n(h) \geq \lambda_0/2$ .

On the other hand in view of (4.7.8)  $\min_{n \leq n_0} \inf_{h \in \mathcal{H}_n} \lambda_n(h) > 0$ . It remains to define  $\lambda > 0$  as

$$\lambda = \min \left( \min_{n \leq n_0} \inf_{h \in \mathcal{H}_n} \lambda_n(h), \lambda_0/2 \right).$$

■

**Proof of Lemma 10** Remind that  $h_k \leq h_\kappa \leq h^*$  by definition of  $h_k$ ,  $h^*$  and  $\kappa$  (see (4.5.23)). Using Proposition 5 with  $h = h_k$ , it yields

$$\begin{aligned} \mathbb{E}_f |\hat{f}^{(k)}(y) - f(y)|^q \mathbb{I}_G &\leq C_q^* \left( \frac{1 \vee Ld n h_k^{\beta+d}}{n h_k^d} \right)^q \\ (4.7.11) \qquad \qquad \qquad &\leq C_q^* \left( \frac{1 \vee Ld n (h^*)^{\beta+d}}{n h_k^d} \right)^q. \end{aligned}$$

The control of  $n(h^*)^{\beta+d}$  requires the following calculation.

$$(4.7.12) \qquad n(h^*)^{\beta+d} \leq 1 + \frac{b - \beta}{(b + d)(\beta + d)} \ln n = \rho_n(\beta)$$

where  $\rho_n(\beta)$  is the price to pay for adaptation defined in (4.1.10). By definition of  $h_k$ , we have

$$\begin{aligned} 1 + \kappa \ln 2 &= 1 + \ln \frac{h_{\max}}{h_k} \geq 1 + \ln \frac{h_{\max}}{h^*} \\ &\geq 1 + \frac{b - \beta}{(b + d)(\beta + d)} \ln n - \frac{1}{\beta + d} \ln [c(1 + (b - \beta) \ln n)]. \end{aligned}$$

Using the classical inequality  $\ln(1+x) \leq x$  and  $c \leq 1/(\beta+d)(b+d)$ , we obtain with the last inequality

$$(4.7.13) \quad \frac{\beta+d-1}{\beta+d} \rho_n(\beta) \leq 1 + \kappa \ln 2 \leq 1 + k \ln 2, \forall k \geq \kappa.$$

According to (4.7.11), (4.7.12) and (4.7.13), Lemma 10 is proved.  $\blacksquare$

**Proof of Lemma 11** Note that for any  $k \geq \kappa + 1$  and by definition of  $\hat{k}$  (4.3.10)

$$\{\hat{k} = k\} = \cup_{l \geq k} \left\{ |\hat{f}^{(k-1)}(y) - \hat{f}^{(l)}(y)| > \hat{M} S_n(l) \right\}.$$

Note that  $S_n(l)$  is monotonically increasing in  $l$  and, therefore,

$$\begin{aligned} \{\hat{k} = k\} &\subseteq \left\{ |\hat{f}^{(k-1)}(y) - f(y)| > 2^{-1} \hat{M} S_n(k-1) \right\} \\ &\cup \left[ \cup_{l \geq k} \left\{ |\hat{f}^{(l)}(y) - f(y)| > 2^{-1} \hat{M} S_n(l) \right\} \right]. \end{aligned}$$

Taking into account that the event  $G$  implies the realization of the event  $\hat{M} \geq M(f)/2 \geq A/2$  we come to the following inequality: for any  $k \geq \kappa + 1$

$$(4.7.14) \quad \begin{aligned} \mathbb{P}(\hat{k} = k, G) &\leq \mathbb{P} \left\{ |\hat{f}^{(k-1)}(y) - \hat{f}(y)| > 4^{-1} M(f) S_n(k-1), G \right\} \\ &+ \sum_{l \geq k} \mathbb{P} \left\{ |\hat{f}^{(l)}(y) - f(y)| > 4^{-1} M(f) S_n(l), G \right\}. \end{aligned}$$

Note that the definition of  $S_n(l)$  yields

$$n h_l^d S_n(l) \geq 432 D_b^3 (32 q d + 16) \lambda^{-1}(h_l) [1 + \ln(h_{\max}/h_l)].$$

Thus, applying Proposition 4 and Lemma 9 with  $J_1 = L d (1 + 6 D_b^2) / 6 A(f) D_b^2$  and  $\varepsilon = J_1 \lambda^{-1}(h_l) [1 + \ln(h_{\max}/h_l)]$ , we obtain  $\forall l \geq k-1$

$$(4.7.15) \quad \begin{aligned} \mathbb{P} \left\{ |\hat{f}^{(l)}(y) - f(y)| > (M(f)/4) S_n(l), G \right\} &\leq \mathfrak{B}(A, M) \mathcal{E}(h_l) [h_{\max}/h_l]^{-8 q d - 4} \\ &= \mathfrak{B}(A, M) \mathcal{E}(h_l) 2^{-l(8 q d + 4)}. \end{aligned}$$

Note that  $b_{h_l} \leq L d h_l^\beta$  since  $f \in \mathbb{H}_d(\beta, L, A, M)$  and, therefore,

$$(4.7.16) \quad \mathcal{N}(h_l) \leq L d n(h_l)^{\beta+d} \leq L d n(h_\kappa)^{\beta+d} \leq L d n(h^*)^{\beta+d}, \quad \forall l \geq k-1.$$

Here we have also used that  $k \geq \kappa + 1$ . We obtain from (4.7.14), (4.7.15) and (4.7.16) that  $k \geq \kappa + 1$

$$\mathbb{P}(\hat{k} = k, G) \leq J_2 \mathfrak{B}(A, M) \exp \left\{ J_1 n(h^*)^{\beta+d} \right\} 2^{-(k-1)(8 q d + 4)},$$

where  $J_2 = (1 - 2^{-(8 q d + 4)})^{-1}$ .  $\blacksquare$

**Proof of Lemma 12** Put for any  $p \in \mathcal{P}_b$

$$W_{ni}^p(y) = p_1! \dots p_d! \frac{h_{\max}^{d-|p|}}{n} K^\top(0) \mathcal{M}_{nh_{\max}}^{-1}(y) K \left( \frac{X_i - y}{h_{\max}} \right) \mathbb{I}_{V_{\max}(y)}(X_i),$$

and note that  $\tilde{\delta}_p = \sum_{i=1}^d 2Y_i W_{ni}^p(y)$ .

The model (4.1.1) can be rewritten as  $2Y_i = f(X_i) + f(X_i)(2U_i - 1)$ . Thus, putting  $F(X) = (f(X_i))_{i=1, \dots, n}$ ,  $V(X) = (f(X_i)(2U_i - 1))_{i=1, \dots, n}$  and

$$\mathcal{D}(f) = \left( \frac{\partial^{|p|} f(y)}{\partial y_1^{p_1} \dots \partial y_d^{p_d}}, p \in \mathcal{P}_\beta \right),$$

we obtain

$$|\hat{M} - M(f)| \leq \|\tilde{\delta} - \mathcal{D}(f)\|_1 \leq \|\mathcal{V} F(X) - \mathcal{D}(f)\|_1 + \|\mathcal{V} V(X)\|_1.$$

Here  $\mathcal{V}$  is  $D_b \times n$ -matrix of general term  $\mathcal{V}_{pi} = W_{ni}^p(y)$  and  $\|\cdot\|_1$  is the  $\ell_1$ -norm. Let us prove that

$$(4.7.17) \quad \mathbb{P}_f \left\{ |\hat{M} - M(f)| > M(f)/2 \right\} \leq \exp \left\{ -\frac{n^{\frac{b}{b+d}}}{8\vartheta_2^2 D_b^2} \right\}.$$

In view of the result proved in Härlde, Hart, Marron, and Tsybakov [1992] and Tsybakov [2008] there exist  $\vartheta_1, \vartheta_2 > 0$  such that

$$\begin{aligned} \|\mathcal{V} F(X) - \mathcal{D}(f)\|_1 &\leq \vartheta_1 h_{\max}^{\beta - \lfloor \beta \rfloor}, \\ \sup_{i,x} |W_{ni}^p(y)| &\leq \frac{\vartheta_2}{nh_{\max}^d}, \quad p \in \mathcal{P}_\beta. \end{aligned}$$

Remind that  $h_{\max} \xrightarrow{n \rightarrow \infty} 0$  and, therefore,  $\exists n_0$  such that  $\vartheta_1 h_{\max}^{\beta - \lfloor \beta \rfloor} \leq M(f)/4$  for any  $n \geq n_0$ . Note that  $n_0$  can be chosen independent on  $f$  since  $M(f)/4 \geq A/4$ . Thus, we get

$$\begin{aligned} &\mathbb{P}_f \left\{ |\hat{M} - M(f)| > M(f)/2 \right\} \\ &\leq \sum_{p \in \mathbb{N}^d: 0 \leq |p| \leq \beta} \mathbb{P}_f \left\{ \left| \sum_{X_i \in [0,1]^d} f(X_i)(2U_i - 1) W_{ni}^p(y) \right| > \frac{M(f)}{4D_b} \right\}. \end{aligned}$$

Noting that  $|f(X_i)(2U_i - 1) W_{ni}^p(y)| \leq M(f) \frac{\vartheta_2}{nh_{\max}^d}$ , applying Hoeffding inequality Boucheron, Bousquet, and Lugosi [2004] and the last inequality, we obtain

$$\begin{aligned} &\sum_{p \in \mathbb{N}^d: 0 \leq |p| \leq \beta} \mathbb{P}_f \left\{ \left| \sum_{X_i \in [0,1]^d} f(X_i)(2U_i - 1) W_{ni}^p(y) \right| > \frac{M(f)}{4D_b} \right\} \\ (4.7.18) \quad &\leq D_b \exp \left\{ -\frac{nh_{\max}^d}{8\vartheta_2^2 D_b^2} \right\} = D_b \exp \left\{ -\frac{n^{\frac{b}{b+d}}}{8\vartheta_2^2 D_b^2} \right\}. \end{aligned}$$



Therefore (4.7.17) is proved. Since  $|f(y) - A(f)| \leq Ldh_{\max}^\beta \leq A(f)/4$  for  $n \geq n_0$  one has

$$\mathbb{P}_f \left\{ |\hat{A} - A(f)| > A(f)/2 \right\} \leq \mathbb{P}_f \left\{ |\hat{M} - M(f)| > A(f)/4 \right\}.$$

Repeating previous calculations we obtain

$$\begin{aligned} \mathbb{P}_f \left\{ |\hat{A} - A(f)| > A(f)/2 \right\} &\leq D_b \exp \left\{ -\frac{[A(f)]^2 n^{\frac{b}{b+d}}}{16[M(f)]^2 \vartheta_2^2 D_b^2} \right\} \\ (4.7.19) \qquad \qquad \qquad &\leq D_b \exp \left\{ -\frac{An^{\frac{b}{b+d}}}{16M\vartheta_2^2 D_b^2} \right\}. \end{aligned}$$

Since  $\mathbb{P}_f(G^c) \leq \mathbb{P}_f(G_{\hat{A}}^c) + \mathbb{P}_f(G_{\hat{M}}^c)$  the assertion of the lemma follows from (4.7.18) and (4.7.19). ■

# Chapter 5

## Huber Estimation

Ce chapitre correspond à l'article [Chichignoud \[2010c\]](#). L'estimateur de Huber est défini pour la première fois dans le cas non-paramétrique. Une caractéristique importante de cet estimateur est sa *robustesse*. En effet, nous démontrons que l'estimateur de Huber estime la fonction de régression dans le modèle de régression additive (défini par (1.2.2)) pour toute densité  $g_\xi$  vérifiant les hypothèses 1. Si les paramètres de régularité  $\beta, L$  sont connus, l'estimateur atteint la vitesse  $\varphi_{n,2}(\beta)$ . Pour l'adaptation, la vitesse atteinte est  $\phi_{n,2}(\beta)$ . Le *critère de Huber* s'appuie sur l'idée de la médiane (norme  $\ell_1$ ) et la moyenne (norme  $\ell_2$ ) au voisinage de 0, ce qui permet d'avoir un comportement plus continu que pour le critère  $\ell_1$ . De plus, celui-ci adopte un comportement *robuste*. Nous accordons une importance toute particulière à la proposition 7 qui est un résultat sur les grandes déviations de l'estimation de Huber.

### 5.1 Introduction

Let statistical experiment is generated by the observation  $Z^{(n)} = (X_i, Y_i)_{i=1, \dots, n}, n \in \mathbb{N}^*$ , where  $(X_i, Y_i)$  satisfies the equation

$$(5.1.1) \quad Y_i = f(X_i) + \xi_i, \quad i = 1, \dots, n.$$

Here  $f : [0, 1]^d \rightarrow \mathbb{R}$  is unknown function and we are interested in estimating  $f$  at a given point  $y \in [0, 1]^d$  from observation  $Z^{(n)}$ .

The random variables (noise)  $(\xi_i)_{i=1, \dots, n}$  are supposed to be independent and identically distributed of law  $\xi$  of unknown density  $g_\xi(\cdot)$  with respect to the Lebesgue's measure  $\mathbb{R}$ .

The design points  $(X_i)_{i=1, \dots, n}$  are random, independent and uniformly distributed on  $[0, 1]^d$ . The random vector  $(X_i)_{i=1, \dots, n}$  is independent of  $(\xi_i)_{i=1, \dots, n}$ .

Along the paper, the unknown function  $f$  is supposed to be smooth, in particular, it belongs to the Holder ball of functions  $\mathbb{H}_d(\beta, L, M)$  (see Definition 16 below). Here  $\beta > 0$

is the smoothness of  $f$ ,  $M$  is the upper bound of  $f$  and its partial derivatives and  $L > 0$  is Lipschitz constant.

The main goal is to estimate  $f$  at a given point  $y \in [0, 1]^d$ , for it, we assume following conditions on the unknown density  $g_\xi$  of the noise.

**Assumptions 2.**

1.  $g_\xi$  is symmetric,
2.  $g_\xi(0) \geq A > 0$ ,
3.  $g_\xi$  is continuous in 0.

Remark that we do not need the existence of the mean for the noise. We suppose in the first assumption that the noise is symmetric. It is reasonable in practice. The other both assumptions imply that it exists a unique theoretical median equals 0 of the noise. i.e. 0 is the unique solution of the equation  $\mathbb{P}(\xi < 0) = 1/2$ .

**Motivation.** In this paper, we deal with robust nonparametric estimation, in particular the idea of the median. We search an estimator which is not sensitive to outliers. It is well-known that the median is more robust than average (see [Rousseeuw and Leroy \[1987\]](#) or [Huber and Ronchetti \[2009\]](#)). In imaging, people use more and more of median filters to reconstruct images ([Astola, Egiazarian, Foi, and Katkovnik \[2010\]](#)).

This idea has also been used by [Arias-Castro and Donoho \[2009\]](#). The robustness of the median is studied for estimation of the non-continuous functions, the authors developed a new method so-called *Two-scale median filter* to improve the minimax rate of convergence.

The asymptotic normality of the median is used by [Brown, Cai, and Zhou \[2008\]](#) to approximate the regression model by the *wavelet sequence data* and use *BlockJS wavelet* (see [Stein \[1981\]](#) et [Cai \[1999\]](#)) for adaptation. But their conditions on the density  $g_\xi$  are stronger than Assumptions 2 and only unidimensional functions are studied. Another new point is studied in this paper : the random design of uniform law on a compact of  $\mathbb{R}^d$ . It is the first time where an estimator, which does not depend on the density  $g_\xi$ , is studied from the initial observations.

The principle of the median falls within the framework of M-estimation [Van de Geer \[2000\]](#). Indeed, it is well-known that the median is the solution to the problem of minimization of absolute values. The main problem is the discontinuity of absolute values of the derivative at 0. In this sense, Huber developed the so-called *Huber function* [Huber \[1964, 1981\]](#) for the estimation. This theory has been widely taken up by [Rousseeuw \[1984\]](#) in parametric models.

Today, many areas of applications use the *Huber function* in relative GPS positioning [Chang and Guo \[2005\]](#) or imaging [Petrus \[1999\]](#). In nonparametric estimation, the Huber approach is introduced by [Tsybakov \[1982a, 1982b, 1983, 1986\]](#) and [Härdle and Tsybakov \[1988, 1992\]](#). [Hall and Jones \[1990\]](#) and [Härdle and Tsybakov \[1992\]](#) used respectively the

*Cross Validation* and *Plug-in method* for the  $L_2$ -risk and pointwise asymptotic normality. Only, [Reiss, Rozenholc, and Cuenod \[2009\]](#) proposed the adaptation for the pointwise risk over locally almost constant functions (i.e.  $\beta \leq 1$ ).

Recently, several approaches of the selection from the family of linear estimators were proposed, see for instance [Goldenshluger and Lepski \[2009a, 2009b\]](#) and [Juditsky, Lepski, and Tsybakov \[2009\]](#) and the references therein. These technologies are used in the construction of data-driven (adaptive) procedures and they are heavily based on the linearity property. As we already mentioned, the *locally Huber estimators* are completely non-linear and in Section 5.3 we propose the selection rule from this family. It requires, in particular, to develop new non-asymptotical exponential inequalities, which may have an independent interest (Proposition 7).

We propose a method of estimation based on the idea of median estimator with the *Huber function* (5.1.6). Actually, we use a locally parametric method, we approximate the function  $f$  by a polynomial  $f_\theta$  on a neighborhood  $V_y(h)$  of size  $h^d$  and we estimate the coefficients  $\theta$  by the *Huber estimator* (5.1.7). We study the maximal risk on the isotropic Hölder spaces and we develop an adaptive procedure based on the Lepski's method for the pointwise estimation.

A possible result of this work is the study of estimating anisotropic functions. Indeed, the methods developed by [Kerkycharian, Lepski, and Picard \[2001\]](#), [Klutchnikoff \[2005\]](#) and [Goldenshluger and Lepski \[2008, 2009a\]](#) are based on the linear properties and they do not agree to these estimators. Another perspective is looking for an oracle inequality for the family of Huber estimators indexed by the bandwidth.

**Huber estimation and maximal risk.** The first part of the paper is devoted to estimate over  $\mathbb{H}_d(\beta, L, M)$  estimation, in particular, when the parameters  $\beta, L, M$  and  $A$  are supposed to be known *a priori*. We find the *rate of convergence* (5.1.2) on  $\mathbb{H}_d(\beta, L, M)$  and propose the estimator with a deterministic choice of  $h_n(\beta, L)$  which achieves the rate (5.1.2). We show that for any  $\beta \in \mathbb{R}_+^*$ , the rate of convergence is bounded from below by the sequence

$$(5.1.2) \quad \varphi_n(\beta) = n^{-\frac{\beta}{2\beta+d}}.$$

We note that (5.1.2) is the rate of linear estimators. But it is not necessarily optimal according to the density of the noise. For example, if  $\xi$  is an uniform random variable, it is well-known that the minimax rate is  $n^{-\frac{\beta}{\beta+d}}$ . In the Gaussian model, we have the same rate  $n^{-\frac{\beta}{2\beta+d}}$  which is minimax.

To construct our estimator, we use so-called *local Huber estimation* which consists in the following. Let

$$V_h(y) = \bigotimes_{j=1}^d [y_j - h/2, y_j + h/2],$$

be the neighborhood around  $y$  such that  $V_h(y) \subseteq [0, 1]^d$ , where  $h \in (0, 1)$  be a given scalar. Fix  $b > 0$  (without loss of generality we will assume that  $b$  is integer) and let

$$(5.1.3) \quad D_b = \sum_{m=0}^b \binom{m+d-1}{d-1}.$$

Let  $K(z), z \in \mathbb{R}^d$  be the  $D_b$ -dimensional vector of polynomials of the following type (the sign  $\top$  below means the transposition):

$$K^\top(z) = \left( \prod_{j=1}^d z_j^{p_j}, \quad (p_1, \dots, p_d) \in \mathbb{N}^d : 0 \leq p_1 + \dots + p_d \leq b \right).$$

For any  $t \in \mathbb{R}^{D_b}$  we define the local polynomial

$$(5.1.4) \quad f_t(x) = t^\top K \left( \frac{x-y}{h} \right) \mathbb{I}_{V_h(y)}(x), \quad x \in [0, 1]^d,$$

where  $\mathbb{I}$  denotes the indicator function. The local polynomial  $f_t$  can be viewed as an approximation of the regression function  $f$  inside of the neighborhood  $V_h$ . Introduce the following subset of  $\mathbb{R}^{D_b}$

$$(5.1.5) \quad \Theta(M) = \{t \in \mathbb{R}^{D_b} : \|t\|_1 \leq M\},$$

where  $\|\cdot\|_1$  is  $\ell_1$ -norm on  $\mathbb{R}^{D_b}$ . Consider the *Huber function*,

$$(5.1.6) \quad Q(z) = \frac{z^2}{2} \mathbb{I}_{|z| \leq 1} + \left( |z| - \frac{1}{2} \right) \mathbb{I}_{|z| > 1},$$

where  $\mathbb{I}$  denotes the indicator function. Put the *Huber criterion*

$$(5.1.7) \quad \tilde{m}_h(t) = \tilde{m}_h(t, Z^{(n)}) := \frac{1}{nh^d} \sum_{i=1}^n Q(Y_i - f_t(X_i)) \mathbb{I}_{\{X_i \in V_h(y)\}}$$

Let  $\check{\theta}(h)$  be the solution of the following minimization problem :

$$(5.1.8) \quad \check{\theta}(h) = \arg \min_{t \in \Theta(M)} \tilde{m}_h(t).$$

The *locally Huber estimator*  $\check{f}^h(y)$  of  $f(y)$  is defined now as  $\check{f}^h(y) = \check{\theta}_{0, \dots, 0}(h)$ . We note that similar locally parametric approach based on maximum likelihood estimators was recently proposed in [Katkovnik and Spokoiny \[2008\]](#) for *regular statistical models* or based on bayesian estimators was recently proposed in [Chichignoud \[2010a\]](#) for *non-linear models*.

As we see, our construction contains an extra-parameter  $h$  to be chosen. To make this choice, we use quite standard arguments. First, we note that in view of  $f \in \mathbb{H}_d(\beta, L, M)$

$$\exists \theta = \theta(f, y, h) \in \Theta(M) : \sup_{x \in V_h(y)} |f(x) - f_\theta(x)| \leq Ldh^\beta.$$

Thus, if  $h$  is chosen sufficiently small our original model (5.1.1) is well approximated inside of  $V_h(y)$  by the “parametric” model

$$\mathcal{Y}_i = f_\theta(X_i) + \xi_i, \quad i = 1, \dots, N_h, \quad N_h = (nh^d)^{1/2},$$

in which the *Huber estimator*  $\check{\theta}$  achieves the rate  $N_h^{-1}$  (See Theorem 15).

Finally,  $h_n(\beta, L)$  is chosen as the solution of the following minimization problem

$$(5.1.9) \quad N_h^{-1} + Ldh^\beta \rightarrow \min_h.$$

and we show that corresponding estimator  $\check{f}^{h_n(\beta, L)}(y)$  achieves (5.1.2) for  $f(y)$  on  $\mathbb{H}_d(\beta, L, M)$ . For the upper bounds (maximal risk and adaptation), we establish majoration of probability:

$$\mathbb{P}_f \left( |\check{f}^h(y) - f(y)| \geq \frac{\varepsilon}{N_h} \right) \leq C_1 \exp \left\{ -C_2 \frac{\varepsilon^2}{A + B\varepsilon} \right\}.$$

With Bernstein inequality, the well known point is  $B\varepsilon$ . For adaptation,  $\varepsilon$  can be large ( $\sqrt{\ln n}$ ), in general, the term  $B = B_n$  is such that  $B_n \sqrt{\ln n} \xrightarrow{n \rightarrow \infty} 0$ .

**Adaptive estimation.** The second part of the paper is devoted to the adaptive minimax estimation over collection of *isotropic* functional classes in the model (5.1.1). To our knowledge, the problem of adaptive estimation, with Huber estimator in pointwise estimation, is not study in the literature. As we mentioned above we consider only the case  $\beta \in (0, b]$  and  $L \in [l_*, l^*]$  will be studied.

Well-known disadvantage of maximal approach is the dependence of the estimator on the parameters describing functional class on which the maximal risk is determined. In particular,  $h_n(\beta, L)$  optimally chosen in view of (5.1.9) depends explicitly on  $\beta$  and  $L$ . To overcome this drawback, the maximal adaptive approach was proposed by Lepski [1990, 1991] and Lepski, Mammen, and Spokoiny [1997]. The first question arising in the adaptation (reduced to the problem at hand) can be formulated as follows.

*Does there exist an estimator which would be minimax on  $\mathbb{H}(\beta, L, M)$  simultaneously for all values of  $\beta$  and  $L$  belonging to some given set  $\mathfrak{B} \subseteq [\mathbb{R}_+ \setminus 0] \times [\mathbb{R}_+ \setminus 0]$  ?*

We can show that the answer of this question is **negative**, that is typical for the estimation of the function at a given point Lepski and Spokoiny [1997]. This answer can be reformulated in the following manner: the family of rates of convergence  $\{\varphi_n(\beta, L), (\beta, L) \in \mathfrak{B}\}$  is **unattainable** for the problem under consideration.

Thus, we need to find another family of normalizations for maximal risk which would be attainable and, moreover, optimal in view of some criterion of optimality. Nowadays, the most developed criterion of optimality is due to [Klutchnikoff \[2005\]](#) in the Gaussian model.

The most important step is to find an estimator, called *adaptive*, which attains this family of normalizations:

$$(5.1.10) \quad \phi_n(\beta, L) = \left( \frac{\rho_n(\beta, L)}{n} \right)^{\frac{\beta}{2\beta+d}}$$

with  $\rho_n(\beta, L) = \left( 1 + \frac{2(b-\beta)}{(2\beta+d)(2b+d)} \ln n \right)^{1/2}$ . The factor  $\rho_n$  can be considered as *price to pay for adaptation* [Lepski \[1990\]](#).

In the present paper, we construct such estimator using general adaptation scheme due to [Lepski, Mammen, and Spokoiny \[1997\]](#) (so-called *Lepski's method*). To our knowledge, it is the first time that this method is applied with *Huber estimator* and in the statistical model with the unknown density of the noise. However, the limitation concerning the consideration of isotropic classes of functions is also due to the use of Lepski's procedure. It seems that to be able to treat the adaptation over the scale of anisotropic classes, another scheme should be applied [Kerkycharian, Lepski, and Picard \[2001\]](#), [Klutchnikoff \[2005\]](#) and [Goldenshluger and Lepski \[2008\]](#). For the model (5.1.1) this problem is still open.

This paper is organized as follows. In Section 5.2 we present the results concerning maximal risk and Section 5.3 is devoted to the adaptive estimation. The proofs of main results are given in Section 5.4 (upper bound), technical lemmas are postponed to Appendix.

## 5.2 Maximal Risk on $\mathbb{H}_d(\beta, L, M)$

In this section, we present several results concerning Maximal risk. We propose the estimator which is bounded the maximal risk on this class of functions under some additional restrictions imposed on these parameters. For any  $(p_1, \dots, p_d) \in \mathbb{N}^d$  we denote  $\vec{p} = (p_1, \dots, p_d)$  and  $|\vec{p}| = p_1 + \dots + p_d$ .

**Definition 16.** Fix  $\beta > 0$ ,  $L > 0$ ,  $M > 0$  and let  $\lfloor \beta \rfloor$  be the largest integer strictly less than  $\beta$ . The isotropic Hölder class  $\mathbb{H}_d(\beta, L, M)$  is the set of functions  $f : [0, 1]^d \rightarrow \mathbb{R}$  having on  $[0, 1]^d$  all partial derivatives of order  $\lfloor \beta \rfloor$  and such that  $\forall x, y \in [0, 1]^d$

$$\sum_{m=0}^{\lfloor \beta \rfloor} \sum_{|\vec{p}|=m} \sup_{x \in [0, 1]^d} \left| \frac{\partial^{|\vec{p}|} f(x)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} \right| \leq M,$$

$$\left| \frac{\partial^{|\vec{p}|} f(x)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} - \frac{\partial^{|\vec{p}|} f(y)}{\partial y_1^{p_1} \dots \partial y_d^{p_d}} \right| \leq L [\|x - y\|_1]^{\beta - \lfloor \beta \rfloor}, \quad \forall |\vec{p}| = \lfloor \beta \rfloor.$$

To measure the performance of estimation procedures on  $\mathbb{H}_d(\beta, L)$  we will use maximal approach.

Let  $\mathbb{E}_f = \mathbb{E}_f^n$  be the mathematical expectation with respect to the probability law of the observation  $Z^{(n)}$  satisfying (5.1.1). Firstly, we define the maximal risk on  $\mathbb{H}_d(\beta, L)$  corresponding to the estimation of the function  $f$  at a given point  $y \in [0, 1]^d$ .

Let  $\tilde{f}_n$  be an arbitrary estimator built from the observation  $Z^{(n)}$ . Set  $\forall r > 0$

$$R_{n,r}[\tilde{f}_n, \mathbb{H}_d(\beta, L)] = \sup_{f \in \mathbb{H}_d(\beta, L)} \mathbb{E}_f |\tilde{f}_n(y) - f(y)|^r.$$

The quantity  $R_{n,r}[\tilde{f}_n, \mathbb{H}_d(\beta, L)]$  is called *maximal risk* of the estimator  $\tilde{f}_n$  on  $\mathbb{H}_d(\beta, L)$ .

**Upper bound for maximal risk.** Let us start with the construction of our estimator based on *local Huber estimation*.

Let

$$(5.2.1) \quad \bar{h} = (L^2 n)^{-\frac{1}{2\beta+d}}, \quad j = 1, \dots, d.$$

The next theorem shows how to construct the estimator basing on locally Huber approach which achieves the rate (5.1.2). Let  $\bar{f}^{\bar{h}}(y) = \check{\theta}_{0,\dots,0}(\bar{h})$  is given by (5.1.5), (5.1.7) and (5.1.8) with  $h = \bar{h}$ .

**Theorem 15.** *Let  $\beta > 0$ ,  $L > 0$  and  $M > 0$  be fixed. Then for any density  $g_\xi$  verified Assumptions 2, we have*

$$\limsup_{n \rightarrow \infty} \varphi_n^{-r}(\beta) R_{n,r}[\bar{f}^{\bar{h}}(y), \mathbb{H}_d(\beta, L, M)] < \infty, \quad \forall r \geq 1,$$

**Remark 15.** *We can see that, asymptotically, the rate of convergence of risk is  $\varphi_n(\beta)$ . We can show in some models that  $n^{-\frac{\beta}{2\beta+d}}$  is minimax (gaussian and Cauchy models, [Tsybakov \[2008\]](#)). [Tsybakov \[1982a\]](#) showed lower bounds for minimax probability risk over Hölder classes with the additive model (5.1.1). But if  $g_\xi$  is the density of uniform random variable, then the minimax rate is  $n^{-\frac{\beta}{\beta+d}}$ , so our Huber estimator is not optimal in this sense. The advantage of this approach is that we do not assume knowledge of the density  $g_\xi$ , only a few conditions (see Assumptions 2), we talk “robust estimation”.*

## 5.3 Bandwidth Selector of Huber Estimator

This section is devoted to the adaptive estimation over the collection of the classes  $\{\mathbb{H}_d(\beta, L, M)\}_{\beta, L}$ , we suppose  $M$  known. We will not impose any restriction on possible



value of  $L$ , but we will assume that  $\beta \in (0, b]$ , where  $b$ , as previously, is an arbitrary a priori chosen integer.

We start with formulating the remark showing that there is not optimally adaptive estimator (here we follow the terminology introduced in Lepski [1990, 1992a, 1992b]). It means that there is not an estimator which would be minimax simultaneously for several values of parameter  $\beta$  even if  $L$  is supposed to be fixed. This result does not require any restriction on  $\beta$  as well. Let  $\Psi = \{\psi_n(\beta)\}_{\beta \in (0, b]}$  be a given family of normalizations.

**Definition 17.** *The family  $\Psi$  is called admissible if there exists an estimator  $\check{f}_n$  such that for some  $L, M > 0$*

$$(5.3.1) \quad \limsup_{n \rightarrow \infty} \psi_n^{-r}(\beta) R_{n,r}(\check{f}_n, \mathbb{H}_d(\beta, L, M)) < \infty, \quad \forall \beta \in (0, b].$$

*The estimator  $\check{f}_n$  satisfying (5.3.1) is called  $\Psi$ -attainable. The estimator  $\check{f}_n$  is called  $\Psi$ -adaptive if (5.3.1) holds for any  $L > 0$  and  $M > 0$ .*

**Remark 16.** *Note that the family of rates of convergence  $\{\varphi_n(\beta)\}_{\beta \in (0, b]}$  is not admissible. We can show this assertion with Lepski and Spokoiny [1997] in Gaussian model.*

Let  $\Phi$  be the following family of normalizations:

$$\phi_n(\beta) = \left( \frac{\rho_n(\beta)}{n} \right)^{\frac{\beta}{2\beta+d}}, \quad \rho_n(\beta) = \left( 1 + \frac{2(b-\beta)}{(2\beta+d)(2b+d)} \ln n \right)^{1/2}, \quad \beta \in (0, b].$$

We remark that  $\phi_n(b) = \varphi_n(b)$  and  $\rho_n(\beta) \sim \sqrt{\ln n}$  for any  $\beta \neq b$ .

**Remark 17.** *We can show that this family  $\Phi$  is adaptive optimal using Klutchnikoff's Criterion Klutchnikoff [2005] in Gaussian model.*

**Construction of  $\Phi$ -adaptive estimator.** As it was already mentioned in Introduction, the construction of our estimation procedure consists of several steps. First, we determine family of *locally Huber estimators*. Next, based on so-called *Lepski's method*, we propose data-driven selection from this family.

We take  $\check{f}^h$  the estimator given by (5.1.5), (5.1.7) and (5.1.8), so the family of locally bayesian estimator  $\check{\mathcal{F}}$  is defined now as follows. Put  $h_{\max} = n^{-\frac{1}{2b+d}}$  and

$$h_k = 2^{-k} h_{\max}, \quad k = 0, \dots, \mathbf{k}_n,$$

where  $\mathbf{k}_n$  is the largest integer such that  $h_{\mathbf{k}_n} \geq h_{\min} = n^{-1/d} \ln^{\frac{b}{d(2b+d)}} n$ . Set

$$(5.3.2) \quad \check{\mathcal{F}} = \left\{ \check{f}^{(k)}(y) = \check{\theta}_{0, \dots, 0}(h_k), \quad k = 0, \dots, \mathbf{k}_n \right\}.$$

We put  $\check{f}^*(y) = \check{f}^{(\hat{k})}(y)$ , where  $\check{f}^{(\hat{k})}(y)$  is selected from  $\check{\mathcal{F}}$  in accordance with the rule:

$$(5.3.3) \quad \hat{k} = \inf \left\{ k = \overline{0, \mathbf{k}_n} : |\check{f}^{(k)}(y) - \check{f}^{(l)}(y)| \leq C S_n(l), \quad l = \overline{k+1, \mathbf{k}_n} \right\}.$$

Here we have used the following notations.

$$(5.3.4) \quad C = 8rd \frac{D_b^2}{\lambda} \left( 8 + \frac{4\sqrt{\lambda}}{3D_b} \right), \quad S_n(l) = \left[ \frac{1 + l \ln 2}{n(h_l)^d} \right]^{1/2}, \quad l = 0, 1, \dots, \mathbf{k}_n,$$

where  $\lambda > 0$  is a constant defined in Lemma 13 and  $D_b$  defined in (5.1.3).

**Theorem 16.** *Let  $b > 0, M > 0$  be fixed. Then, for any density  $g_\xi$  verified Assumptions 2, for any  $\beta \in (0, b]$ ,  $L > 0$  and  $r \geq 1$*

$$\limsup_{n \rightarrow \infty} \phi_n^{-r}(\beta) R_{n,r} \left[ \check{f}^*(y), \mathbb{H}_d(\beta, L, M) \right] < \infty.$$

**Remark 18.** *The assertion of the theorem means that the proposed estimator  $\check{f}^*(y)$  is  $\Phi$ -adaptive. It implies in particular that the family of normalizations  $\Phi$  is admissible.*

**Remark 19.** *Remark that  $M$  is involved in construction of the set of coefficients  $\Theta(M)$ , for example we can replace  $M$  by  $\ln n$ . The constant  $\lambda$  depends on  $g_\xi(0)$  and  $\|g_\xi\|_\infty$  (see Proof of Lemma 13), unknown in practice, but we can calibrate the constant  $C$ .*

## 5.4 Proofs of Main Results: Upper Bounds

Let  $\mathcal{H}_n, n > 1$  be the following subinterval of  $(0, 1)$ .

$$(5.4.1) \quad \mathcal{H}_n = \left[ \frac{(\ln n)^{\frac{1}{d}}}{n^{1/d}}, \left( \frac{1}{\ln n} \right)^{\frac{1}{2b+d}} \right].$$

Later on, we will consider only the values of  $h$  belonging to  $\mathcal{H}_n$ . Put also  $N_h = (nh^d)^{1/2}$ .

### 5.4.1 Auxiliary Results: Large Deviations for M-estimators

Let  $h \in \mathcal{H}_n$ , put  $\theta = \theta(f, y, h) = \{\theta_{\vec{p}} : \vec{p} \in \mathbb{N}^d : 0 \leq |\vec{p}| \leq b\}$ , where  $\theta_0 = \theta_{0,\dots,0} = f(y)$  and

$$\theta_{\vec{p}} = \frac{\partial^{|\vec{p}|} f(y)}{\partial y_1^{p_1} \dots \partial y_d^{p_d}} \frac{h^{|\vec{p}|}}{p_1! \dots p_d!}, \quad p \in \mathbb{N}^d : 0 < |\vec{p}| \leq b.$$

Remind the agreement which we follow in the present paper: if the function  $f$  and vector  $\vec{p}$  be such that  $\partial^{|\vec{p}|} f$  does not exist we put  $\theta_{\vec{p}} = 0$ .

Let  $f_\theta(x)$ , given by (5.1.4), be the local polynomial approximation of  $f$  inside  $V_h(y)$  and let  $b_h$  be the corresponding approximation error, i.e.

$$b_h := \sup_{x \in V_h(y)} |f_\theta(x) - f(x)|.$$

Introduce for any function  $f$

$$(5.4.2) \quad M(f) = \sum_{m=0}^b \sum_{p_1+\dots+p_d=m} \left| \frac{\partial^m f(y)}{\partial x_1^{p_1} \dots \partial x_d^{p_d}} \right|.$$

The following agreement will be used in the sequel: if the function  $f$  and  $m \geq 1$  be such that  $\partial^m f$  does not exist, we will put formally  $\partial^m f = 0$  in the definition of  $M(f)$ .

Put

$$(5.4.3) \quad q(z) := Q'(z) = \mathbb{I}_{]1,+\infty[}(z) - \mathbb{I}_{]-\infty,-1[}(z) + z \mathbb{I}_{[-1,1]}(z), \quad z \in \mathbb{R},$$

where  $Q'$  is the derivative of *Huber function*  $Q$  defined in (5.1.6). Set  $\mathcal{S}_b = \{\vec{q} \in \mathbb{N}^d : |\vec{q}| \leq b\}$  and for any  $t \in \Theta(M)$  define the so-called *partial derivative of Huber criterion*

$$(5.4.4) \quad \tilde{D}_h^{\vec{p}}(t) = \frac{\partial^{|\vec{p}|} \tilde{m}_h(t)}{\partial t_{\vec{p}}} = -\frac{1}{nh^d} \sum_{i=1}^n q(Y_i - f_t(X_i)) \left( \frac{X_i - y}{h} \right)^{\vec{p}} \mathbb{I}_{X_i \in V_y(h)}, \quad \vec{p} \in \mathcal{S}_b,$$

where  $x^{\vec{p}} = x_1^{p_1} \dots x_d^{p_d}$ ,  $\forall x \in [0, 1]^d$ . Put the *expectation of partial derivative of Huber criterion*:

$$(5.4.5) \quad \begin{aligned} \mathcal{E}_h^{\vec{p}}(t) &:= \mathbb{E}_f \left[ \tilde{D}_h^{\vec{p}}(t) \right] \\ &= -\frac{1}{h^d} \int_{V_y(h)} \left( \frac{x-y}{h} \right)^{\vec{p}} \left[ \int_{1+f_t(x)-f(x)}^{+\infty} g_\xi(z) dz - \int_{-\infty}^{-1+f_t(x)-f(x)} g_\xi(z) dz \right. \\ &\quad \left. + \int_{-1+f_t(x)-f(x)}^{1+f_t(x)-f(x)} (z - f_t(x) + f(x)) g_\xi(z) dz \right] dx, \quad \vec{p} \in \mathcal{S}_b, \end{aligned}$$

We decompose  $\mathcal{E}_h^{\vec{p}}$  in two parts:  $\mathcal{E}_h^{\vec{p}}(t) = D_h^{\vec{p}}(t) + \Delta_h^{\vec{p}}(t)$ , where

$$(5.4.6) \quad \begin{aligned} D_h^{\vec{p}}(t) &= -\frac{1}{h^d} \int_{V_y(h)} \left( \frac{x-y}{h} \right)^{\vec{p}} \left[ \int_1^{+\infty} g_\xi(z + f_{t-\theta}(x)) dz - \int_{-\infty}^{-1} g_\xi(z + f_{t-\theta}(x)) dz \right. \\ &\quad \left. + \int_{-1}^1 z g_\xi(z + f_{t-\theta}(x)) dz \right] dx, \quad \vec{p} \in \mathcal{S}_b. \end{aligned}$$

and to simplify we note

$$(5.4.7) \quad \Delta_h^{\vec{p}}(t) = \mathcal{E}_h^{\vec{p}}(t) - D_h^{\vec{p}}(t), \quad \vec{p} \in \mathcal{S}_b,$$

$\Delta_h(t) = (\Delta_h^{\vec{p}}(t))_{\vec{p} \in \mathcal{S}_b}^\top$  and  $D_h(t) = (D_h^{\vec{p}}(t))_{\vec{p} \in \mathcal{S}_b}^\top$  are respectively so-called *bias term* and *deterministic criterion*. Finally, let us define the *partial derivative of deterministic criterion* by the following notation

$$(5.4.8) \quad \nabla_h^{\vec{p}, \vec{q}}(t) = \frac{\partial^{|\vec{q}|} D_h^{\vec{p}}(t)}{\partial t_{\vec{q}}} = \frac{1}{h^d} \int_{V_h(y)} \left( \frac{x-y}{h} \right)^{\vec{p}+\vec{q}} \int_{-1}^1 g_\xi(z + f_{t-\theta}(x)) dz dx, \quad \vec{p}, \vec{q} \in \mathcal{S}_b^2.$$

We call *Jacobian matrix of deterministic criterion*

$$(5.4.9) \quad J_D(t) = \left( \nabla_h^{\vec{p}, \vec{q}}(t) \right)_{\vec{p}, \vec{q} \in \mathcal{S}_b^2}.$$

The main property of the matrix  $J_D(\cdot)$  is that it exists a neighborhood  $\mathcal{V}$  of  $\theta$  does not depend on  $n$  such that  $J_D(\cdot)$  is invertible on  $\mathcal{V}$  (see Lemma 13). The next proposition is the milestone for all results proved in the paper.

Put the random sets for any  $h, z > 0$

$$(5.4.10) \quad G_z^h = \{\check{\theta}(h) \in \mathcal{B}(\theta, z)\}, \quad \mathcal{B}(\theta, z) = \{t \in \theta(M) : \|t - \theta\|_2 \leq z\}.$$

**Proposition 7.** *For any  $n \in \mathbb{N}^*$ ,  $h \in \mathcal{H}_n$ ,  $f$  such that  $M(f) < \infty$  and  $\delta > 0$  such that  $\mathcal{B}(\theta, \delta) \subseteq \mathcal{V}$ , then  $\forall \varepsilon \geq D_b \varepsilon_0(h)/\sqrt{\lambda}$  we have*

$$\mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon, G_\delta^h \right) \leq \Sigma \exp \left\{ - \frac{(\varepsilon \sqrt{\lambda} D_b^{-1} - \varepsilon_0(h))^2}{8 + \frac{8}{3N_h} \varepsilon \sqrt{\lambda} D_b^{-1}} \right\},$$

where  $\varepsilon_0(h) = (1 \vee b_h N_h)$ ,  $\check{f}^h(y)$  defined in (5.1.5), (5.1.7), (5.1.8),  $\lambda$  is a positive constant defined in Lemma 13 and  $\Sigma$  is defined in Lemma 15.

The next proposition provides us with upper bound for the risk of a locally Huber estimator. Put

$$(5.4.11) \quad \bar{C}_r = \frac{D_b^r}{\lambda^{r/2}} \left[ 1 + \int_0^{+\infty} (D_b^{-1} \sqrt{\lambda} \eta + 1)^{r-1} \exp \left\{ - \frac{\eta^2}{8(1 + \eta + L \sqrt{\lambda} D_b^{-1})} \right\} d\eta \right], \quad r \geq 1.$$

**Proposition 8.** *For any  $n \in \mathbb{N}^*$ ,  $h \in \mathcal{H}_n$ ,  $\delta > 0$  such that  $\mathcal{B}(\theta, \delta) \subseteq \mathcal{V}$  and  $f$  such that  $M(f) < \infty$ , then*

$$\mathbb{E}_f |\check{f}^h(y) - f(y)|^r \mathbb{I}_{G_\delta^h} \leq \bar{C}_r \varepsilon_0^r(h) N_h^{-r}, \quad r \geq 1.$$

The proof of this proposition is elementary by integration of Proposition 7.

### 5.4.2 Proof of Proposition 7

**Auxiliary lemmas.** We give the following lemma concerning *deterministic criterion*.

**Lemma 13.** *For any  $g_\xi$  verified Assumptions 2, then*

1.  $J_D(\theta)$  is positive definite matrix.
2.  $\theta$  is unique solution of  $D_h(\cdot) = (0, \dots, 0)$ .
3.  $\forall h \in \mathcal{H}_n$ , it exists a neighborhood  $\mathcal{V}$  of  $\theta$ , which does not depend on  $n$ , such that

$$\|\tilde{\theta} - \theta\|_2 \leq \lambda^{-1/2} \|D_h(\tilde{\theta}) - D_h(\theta)\|_2, \quad \tilde{\theta} \in \mathcal{V},$$

where  $\lambda = \inf_{\tilde{\theta} \in \mathcal{V}} \lambda_{\tilde{\theta}}$ ,  $\lambda_{\tilde{\theta}}$  is the smallest eigenvalue of matrix  $J_D(\tilde{\theta})^\top J_D(\tilde{\theta})$  and  $\|\cdot\|_2$  is the  $\ell_2$ -norm on  $\mathbb{R}^{D_b}$ .

The proof is given in Appendix. Let us give a lemma which allows to control the *bias term*.

**Lemma 14.** *For any  $h \in \mathcal{H}_n$  and  $f$  such that  $M(f) < \infty$ , we have*

$$|\Delta_h^{\vec{p}}(t)| \leq b_h, \quad \vec{p} \in \mathcal{S}_b, \quad t \in \theta_M,$$

where  $\|\cdot\|_\infty$  is the sup-norm

The next result allows to control the large deviations of *partial derivatives of Huber criterion*.

**Lemma 15.** *For any  $f$  such that  $M(f) < \infty$  and  $h \in \mathcal{H}_n$ , we have  $\forall z \geq (1 \vee b_h N_h)$*

$$\mathbb{P}_f \left( N_h \sup_{t \in \Theta(M)} |\tilde{D}_h^{\vec{p}}(t) - D_h^{\vec{p}}(t)| \geq z \right) \leq \Sigma \exp \left\{ -\frac{(z - b_h N_h)^2}{8 + \frac{8}{3N_h} z} \right\}, \quad \vec{p} \in \mathcal{S}_b,$$

where  $\Sigma$  is defined in the proof.

The proof (given in Appendix) is based on a chaining argument and Bernstein inequality. It is required that the *Huber criterion* is  $C^1$ .

**Lemma 16.** *For any  $f \in \mathbb{H}_d(\beta, L, M)$ ,  $\delta > 0$  such that  $\mathcal{B}(\theta, \delta) \subseteq \mathcal{V}$  and  $n \in \mathbb{N}^*$ ,  $h \in \mathcal{H}_n$  such that*

$$b_h \leq \frac{\varkappa_\delta}{2\sqrt{D_b}}, \quad N_h \geq \frac{2\sqrt{D_b}}{\varkappa_\delta}, \quad \varkappa_\delta = \inf_{t \in \Theta(M) \setminus \mathcal{B}(\theta, \delta)} \frac{\|D_h(t)\|_2}{2},$$

then

$$\mathbb{P}_f [(G_\delta^h)^c] \leq D_b \Sigma \exp \left\{ -\frac{nh^d \varkappa_\delta^2}{4D_b (8 + 4\varkappa_\delta / (3\sqrt{D_b}))} \right\},$$

where  $(G_\delta^h)^c$  is the additional event of  $G_\delta^h$ ,  $\mathcal{V}$  defined in Lemma 13.

**Proof of Proposition 7.** Remark that  $\tilde{D}_h(\check{\theta}) = (0, \dots, 0)$  where  $\tilde{D}_h(\cdot) = (\tilde{D}_h^{\vec{p}}(\cdot))^{\top}_{\vec{p} \in \mathcal{S}_b}$  defined in (5.4.4). The main idea of this proof is that if  $\tilde{D}_h(t) \xrightarrow[n \rightarrow \infty]{N_h^{-1}} \mathcal{E}_h(t)$  in probability uniformly in  $t$  then we have  $\check{\theta}(h) \xrightarrow[n \rightarrow \infty]{N_h^{-1}} \theta$ . This is the idea of M-estimation [Van de Geer \[2000\]](#). The definition of  $\check{\theta}(h)$  and  $\theta = \theta(f, y, h)$  implies that  $\forall \varepsilon \geq D_b \sqrt{\lambda} \varepsilon_0(h)$

$$\begin{aligned} \mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon, G_{\delta}^h \right) &\leq \mathbb{P}_f \left( N_h |\check{\theta}_0(h) - \theta_0| \geq \varepsilon, G_{\delta}^h \right) \\ &\leq \mathbb{P}_f \left( N_h \sqrt{D_b} \|\check{\theta}(h) - \theta\|_2 \geq \varepsilon, G_{\delta}^h \right). \end{aligned}$$

Under the event  $G_{\delta}^h$  we have  $\check{\theta}(h) \in \mathcal{V}$ , according Lemma 13 we obtain that

$$\mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon, G_{\delta}^h \right) \leq \mathbb{P}_f \left( N_h \frac{\sqrt{D_b}}{\sqrt{\lambda}} \|D_h(\check{\theta}(h)) - D_h(\theta)\|_2 \geq \varepsilon \right).$$

Remember that  $D_h(\theta) = 0$  and  $\tilde{D}_h(\check{\theta}(h)) = 0$ . Using the last inequality, we get

$$\begin{aligned} \mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon, G_{\delta}^h \right) &\leq \mathbb{P}_f \left( N_h \frac{\sqrt{D_b}}{\sqrt{\lambda}} \|D_h(\check{\theta}(h)) - \tilde{D}_h(\check{\theta}(h))\|_2 \geq \varepsilon \right) \\ &\leq \sum_{\vec{p} \in \mathcal{S}_b} \mathbb{P}_f \left( N_h \sup_{t \in \Theta(M)} \left| \tilde{D}_h^{\vec{p}}(t) - D_h^{\vec{p}}(t) \right| \geq \frac{\varepsilon \sqrt{\lambda}}{D_b} \right). \end{aligned}$$

Applying Lemma 15 with  $z = \varepsilon \sqrt{\lambda} / D_b$  and last inequality, finally we have the assertion of Proposition 7

$$\mathbb{P}_f \left( N_h |\check{f}^h(y) - f(y)| \geq \varepsilon, G_{\delta}^h \right) \leq \Sigma \exp \left\{ - \frac{(\varepsilon \sqrt{\lambda} D_b^{-1} - \varepsilon_0(h))^2}{8 + \frac{8}{3N_h} \varepsilon \sqrt{\lambda} D_b^{-1}} \right\}.$$

■

### 5.4.3 Proof of Theorem 15

By definition of  $\bar{h} = (L^2 n)^{-\frac{1}{2\beta+d}}$ , we have  $b_{\bar{h}} N_{\bar{h}} = d$ ,  $N_{\bar{h}}^{-r} = L^{\frac{rd}{2\beta+d}} \varphi_n^r(\beta)$ . So we obtain that the term  $\varepsilon_0(\bar{h})$  defined in Proposition 7 can be written as  $\varepsilon_0(\bar{h}) = d$ . We get

$$(5.4.12) \quad \mathbb{E}_f |\bar{f}^{\bar{h}}(y) - f(y)|^r = \mathbb{E}_f |\bar{f}^{\bar{h}}(y) - f(y)|^r \mathbb{I}_{G_{\delta}^{\bar{h}}} + \mathbb{E}_f |\bar{f}^{\bar{h}}(y) - f(y)|^r \mathbb{I}_{(G_{\delta}^{\bar{h}})^c}.$$

The right hand side is controlled by Lemma 16. Indeed, we can use the Cauchy-Schwarz inequality,

$$\begin{aligned}
 \mathbb{E}_f |\bar{f}^h(y) - f(y)|^r \mathbb{I}_{(G_\delta^h)^c} &\leq \left( \mathbb{E}_f |\bar{f}^h(y) - f(y)|^{2r} \mathbb{P}_f \left\{ (G_\delta^h)^c \right\} \right)^{1/2} \\
 (5.4.13) \qquad \qquad \qquad &\leq (2M)^r D_b \Sigma \exp \left\{ -\frac{nh^d \varkappa_\delta^2}{4D_b (8 + 4\varkappa_\delta/(3\sqrt{D_b}))} \right\}.
 \end{aligned}$$

Using Proposition 8, (5.4.12) and (5.4.13), we obtain

$$\mathbb{E}_f |\bar{f}^h(y) - f(y)|^r \mathbb{I}_{G_\delta^h} \leq \bar{C}_r d L^{\frac{rd}{2\beta+d}} \varphi_n^r(\beta) + (2M)^r D_b \Sigma \exp \left\{ -\frac{nh^d \varkappa_\delta^2}{4D_b (8 + 4\varkappa_\delta/(3\sqrt{D_b}))} \right\}.$$

By passing to the limit in  $n$ , the theorem 15 is proved.  $\blacksquare$

#### 5.4.4 Proof of Theorem 16

We start the proof with formulating some auxiliary results whose proofs are given in Appendix. Define  $J_1 = \sqrt{\lambda}/D_b$  and

$$h^* = \left[ \frac{c \rho_n^2(\beta)}{L^2 d^2 n} \right]^{\frac{1}{2\beta+d}}, \quad c < (2J_1)^{-2}.$$

and let the integer  $\kappa$  be defined as follows.

$$(5.4.14) \qquad \qquad \qquad 2^{-\kappa} h_{\max} \leq h^* < 2^{-\kappa+1} h_{\max}.$$

**Lemma 17.** *For any  $f \in \mathbb{H}_d(\beta, L, M, A)$  and any  $k \geq \kappa + 1$*

$$\mathbb{P}_f(\hat{k} = k, G_\delta^{h_k}) \leq J_2 2^{-2(k-1)rd},$$

where  $J_2 = D_b \Sigma (1 + (1 - 2^{-2rd})^{-1})$ .

**Proof of Theorem 16.** This proof is based on the scheme due to Lepski, Mammen, and Spokoiny [1997]. The definition of  $h^*$  and  $\kappa$  (5.4.14) implies that

$$\varepsilon_0(h_k) \leq D_b \frac{\sqrt{1 + k \ln 2}}{\sqrt{\lambda}}, \quad \forall k \geq \kappa,$$

where  $\varepsilon_0(\cdot)$  defined in Proposition 7. Last inequality yields

$$(5.4.15) \qquad \mathbb{E}_f |\check{f}^{(k)}(y) - f(y)|^r \mathbb{I}_{G_\delta^{h_k}} \leq \bar{C}_r D_b^r \lambda^{-r/2} S_n^r(k), \quad \forall k \geq \kappa.$$

To get this result we have applied Proposition 8 with  $h = h_k$ . We also have

$$\begin{aligned}
& \mathbb{E}_f |f^{(\hat{k})}(y) - f(y)|^r \\
&= \mathbb{E}_f |f^{(\hat{k})}(y) - f(y)|^r \mathbb{I}_{\hat{k} \leq \kappa, G_\delta^{h_{\hat{k}}}} + \mathbb{E}_f |f^{(\hat{k})}(y) - f(y)|^r \mathbb{I}_{\hat{k} > \kappa, G_\delta^{h_{\hat{k}}}} \\
&\quad + \mathbb{E}_f |f^{(\hat{k})}(y) - f(y)|^r \mathbb{I}_{(G_\delta^{h_{\hat{k}}})^c} \\
(5.4.16) \quad &:= R_1(f) + R_2(f) + R_3(f).
\end{aligned}$$

First we control  $R_1$ . Obviously

$$|f^{(\hat{k})}(y) - f(y)|^r \leq 2^{r-1} |f^{(\hat{k})}(y) - f^{(\kappa)}(y)|^r + 2^{r-1} |f^{(\kappa)}(y) - f(y)|^r.$$

The definition of  $\hat{k}$  yields

$$|f^{(\hat{k})}(y) - f^{(\kappa)}(y)|^r \mathbb{I}_{\hat{k} \leq \kappa, G_\delta^{h_{\hat{k}}}} \leq C^r S_n^r(\kappa),$$

In view of (5.4.15) and definitions of  $h_\kappa$  and  $h^*$ , we also get

$$\mathbb{E}_f |f^{(\kappa)}(y) - f(y)|^r \mathbb{I}_{\hat{k} \leq \kappa, G_\delta^{h_{\hat{k}}}} \leq \bar{C}_r D_b^r \lambda^{-r/2} S_n^r(\kappa).$$

Noting that the right hand side of the obtain inequality is independent of  $f$  and taking into account the definition of  $\kappa$  and  $h^*$  we obtain

$$(5.4.17) \quad \limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-r}(\beta) R_1(f) < \infty.$$

Now, let us bounded from above  $R_2$ . Applying Cauchy-Schwarz inequality we have in view of Lemma 17

$$\begin{aligned}
R_2(f) &= \sum_{k=\kappa}^{\mathbf{k}_n} \mathbb{E}_f |f^{(k)}(y) - f(y)|^r \mathbb{I}_{G_\delta^{h_k}} \\
&\leq \sum_{k>\kappa} \left( \mathbb{E}_f |f^{(k)}(y) - f(y)|^{2r} \mathbb{I}_{G_\delta^{h_k}} \right)^{1/2} \sqrt{\mathbb{P}_f \{ \hat{k} = k, G_\delta^{h_k} \}} \\
(5.4.18) \quad &= \sqrt{J_2} \sum_{k>\kappa} \left( \mathbb{E}_f |f^{(k)}(y) - f(y)|^{2r} \mathbb{I}_{G_\delta^{h_k}} \right)^{1/2} 2^{-(k-1)rd}.
\end{aligned}$$

We obtain from (5.4.15) and (5.4.18)

$$(5.4.19) \quad R_2(f) \leq N_{h_{\max}}^{-r} \times D_b \frac{\bar{C}_{2r}^r 2^{rd} \sqrt{J_2}}{(\sqrt{\lambda})^r} \sum_{s \geq 0} (1 + k \ln 2)^{r/2} 2^{-srd}.$$



It remains to note that the right hand side of (5.4.19) is independent of  $f$ . Thus, we have

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-r}(\beta) R_2(f) < \infty.$$

that yields together with (5.4.16) and (5.4.17)

$$\limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-r}(\beta) \mathbb{E}_f \left| \check{f}^{(\hat{k})}(y) - f(y) \right|^r \mathbb{I}_{G_\delta^{h_{\hat{k}}}} < \infty.$$

To get the assertion of the theorem it suffices to show that

$$(5.4.20) \quad \limsup_{n \rightarrow \infty} \sup_{f \in \mathbb{H}_d(\beta, L, A, M)} \phi_n^{-r}(\beta) \mathbb{E}_f \left| \check{f}^{(\hat{k})}(y) - f(y) \right|^r \mathbb{I}_{(G_\delta^{h_{\hat{k}}})^c} < \infty.$$

Note that  $|\check{f}^{(\hat{k})}(y)| \leq M$  by definition. This allows us to state that  $|\check{f}^{(\hat{k})}(y) - f(y)| \leq 2M$ . Here we also have taken into account that  $\|f\|_\infty \leq M$ .

Finally we obtain

$$R_3(f) \leq 2^r M^r \mathbb{P}_f \left\{ (G_\delta^{h_{\hat{k}}})^c \right\}.$$

and (5.4.20) follows now from Lemma 16. ■

## 5.5 Appendix

**Proof of lemma 13.** 1) Remember that for any  $\vec{p}, \vec{q} \in \mathcal{S}_b^2$

$$\nabla_h^{\vec{p}, \vec{q}}(\theta) = \frac{1}{h^d} \int_{V_h(y)} \left( \frac{x-y}{h} \right)^{\vec{p}+\vec{q}} \int_{-1}^1 g_\xi(z) dz dx = \int_{-1}^1 g_\xi(z) dz \int_{[-\frac{1}{2}, \frac{1}{2}]^d} x^{\vec{p}+\vec{q}} dx.$$

So we can write that

$$J_D(\theta) = \int_{-1}^1 g_\xi(z) dz \int_{[-\frac{1}{2}, \frac{1}{2}]^d} K(x) K^\top(x) dx.$$

According to *Assumptions 2*, we have that  $\int_{-1}^1 g_\xi(z) dz > 0$ . Now we show that

$\int_{[-\frac{1}{2}, \frac{1}{2}]^d} K(x) K^\top(x) dx$  is symmetric definite matrix, indeed for any  $\tau \in \mathbb{R}^{D_b} \setminus 0$

$$\tau^\top \int_{[-\frac{1}{2}, \frac{1}{2}]^d} K(x) K^\top(x) dx \tau = \int_{[-\frac{1}{2}, \frac{1}{2}]^d} [\tau^\top K(x)]^2 dx > 0.$$

As  $q'$  is almost continuous everywhere and by definition of  $J_D(\cdot)$  in (5.4.9), it exists a neighborhood of  $\theta$  noted  $\mathcal{V}$  such that the matrix  $J_D(\tilde{\theta})$  is invertible for all  $\tilde{\theta} \in \mathcal{V}$ .

2) With *Assumptions 2*, in particular that  $g_\xi$  is even, by definition of  $D_h$  in (5.4.6) we have that  $D_h(\theta) = 0$ . According to 1), the function  $D_h$  is locally invertible around  $\theta$ , so  $D_h(\theta) = 0$  is unique solution on  $\mathcal{V}$ . Then, we can see that each partial derivative (5.4.8) is positive, thus  $\theta$  is the unique solution of  $D_h(\cdot) = (0, \dots, 0)$  on  $\mathbb{R}^{D_b}$ .

3) We can write the following equalities, for any  $\tilde{\theta} \in \mathcal{V}$  then  $\exists \bar{\theta} \in \mathcal{V}$  such that  $D_h(\tilde{\theta}) - D_h(\theta) = J_D(\bar{\theta})(\tilde{\theta} - \theta)$  and

$$\|D_h(\tilde{\theta}) - D_h(\theta)\|_2^2 = (\tilde{\theta} - \theta)^\top J_D(\bar{\theta})^\top J_D(\bar{\theta})(\tilde{\theta} - \theta).$$

It is elementary to show that  $J_D(\bar{\theta})^\top J_D(\bar{\theta})$  is definite positive matrix if  $J_D(\bar{\theta})$  is invertible. Indeed, for  $\tau \in \mathbb{R}^{D_b} \setminus 0$  such that

$$\tau^\top J_D(\bar{\theta})^\top J_D(\bar{\theta}) \tau = 0 \Rightarrow J_D(\bar{\theta}) \tau = 0 \Rightarrow \tau = 0,$$

it is a contradiction. Thus with 1)

$$(5.5.1) \quad \|D_h(\tilde{\theta}) - D_h(\theta)\|_2 \geq \sqrt{\lambda} \|\tilde{\theta} - \theta\|_2, \quad \tilde{\theta} \in \mathcal{V},$$

where  $\lambda = \inf_{\bar{\theta} \in \mathcal{V}} \lambda_{\bar{\theta}}$  and  $\lambda_{\bar{\theta}}$  is the smallest eigenvalue of matrix  $J_D(\bar{\theta})^\top J_D(\bar{\theta})$ . It follows that  $\lambda > 0$ . ■

**Proof of lemma 14.** By definition of  $D_h^{\vec{p}}$  (5.4.6) and  $\mathcal{E}_h^{\vec{p}}$  (5.4.5),  $\Delta_h^{\vec{p}}$  (5.4.7) can be written as

$$\begin{aligned} \Delta_h^{\vec{p}}(t) = & -\frac{1}{h^d} \int_{V_y(h)} \left( \frac{x-y}{h} \right)^{\vec{p}} \left[ \int_{1+f_t(x)-f(x)}^{1+f_{t-\theta}(x)} g_\xi(z) dz - \int_{-1+f_{t-\theta}(x)}^{-1+f_t(x)-f(x)} g_\xi(z) dz \right. \\ & + \int_{1+f_{t-\theta}(x)}^{1+f_t(x)-f(x)} (z - f_t(x) + f(x)) g_\xi(z) dz \\ & + \int_{-1+f_t(x)-f(x)}^{-1+f_{t-\theta}(x)} (z - f_t(x) + f(x)) g_\xi(z) dz \\ & \left. + \int_{-1+f_{t-\theta}(x)}^{1+f_{t-\theta}(x)} (f(x) - f_\theta(x)) g_\xi(z) dz \right] dx. \end{aligned}$$

With simple calculations, we bound the *bias term*

$$|\Delta_h^{\vec{p}}(t)| \leq b_h.$$

■

**Proof of lemma 15.** First, let us bound the following term for any  $\vec{p} \in \mathcal{S}_b$

$$\sup_{t \in \Theta(M)} \left| \tilde{D}_h^{\vec{p}}(t) - D_h^{\vec{p}}(t) \right| \leq \sup_{t \in \Theta(M)} \left| \tilde{D}_h^{\vec{p}}(t) - \mathcal{E}_h^{\vec{p}}(t) \right| + \sup_{t \in \Theta(M)} \left| \Delta_h^{\vec{p}}(t) \right|,$$

where  $\mathcal{E}_h^{\vec{p}}$  and  $\Delta_h^{\vec{p}}$  are respectively defined in (5.4.5) and (5.4.7). In view of Lemma 14, we get

$$(5.5.2) \quad \sup_{t \in \Theta(M)} \left| \tilde{D}_h^{\vec{p}}(t) - D_h^{\vec{p}}(t) \right| \leq \sup_{t \in \Theta(M)} \left| \tilde{D}_h^{\vec{p}}(t) - \mathcal{E}_h^{\vec{p}}(t) \right| + b_h.$$

Note the function  $L(\cdot) = \tilde{D}_h^{\vec{p}}(\cdot) - \mathcal{E}_h^{\vec{p}}(\cdot)$ . To show the assertion of the lemma, we use a chaining argument on  $L(\cdot)$ . Remember that  $\Theta(M)$  is a compact of  $\mathbb{R}^{D_b}$ . Let  $t_0 \in \Theta(M)$  fixed and for any  $l \in \mathbb{N}^*$  put  $\Gamma_l$  a  $a^{-l}$ -net on  $\Theta(M)$  where the constant  $a = \exp \left\{ \frac{4\pi^{16}}{e^{1842}} \right\}$ . Introduce the following notations

$$u_0(t) = t_0, \quad u_l(t) = \arg \inf_{u \in \Gamma_l} \|u - t\|_1, \quad l \in \mathbb{N}^*.$$

The chaining argument rests on the following decomposition:

$$(5.5.3) \quad L(t) = L(t_0) + \sum_{l=1}^{\infty} L(u_l(t)) - L(u_{l-1}(t)), \quad \forall t \in \Theta(M).$$

Using (5.5.2) and (5.5.3), we obtain

$$(5.5.4) \quad \begin{aligned} & \mathbb{P}_f \left( N_h \sup_{t \in \Theta(M)} |L(t)| \geq z - b_h N_h \right) \\ & \leq \mathbb{P}_f \left( |L(t_0)| + \sup_{t \in \Theta(M)} \sum_{l=1}^{\infty} |L(u_l(t)) - L(u_{l-1}(t))| \geq \frac{z}{N_h} - b_h \right). \end{aligned}$$

We can control the second term as follows.

$$\sup_{t \in \Theta(M)} \sum_{l=1}^{\infty} |L(u_l(t)) - L(u_{l-1}(t))| \leq \sum_{l=1}^{\infty} \sup_{\substack{u, v \in \Gamma_l \times \Gamma_{l-1} \\ \|u-v\|_1 \leq a^{-l}}} |L(u) - L(v)|,$$

where  $\Gamma_0 = \{t_0\}$ . Using (5.5.4) and last inequality, we get

$$(5.5.5) \quad \begin{aligned} & \mathbb{P}_f \left( N_h \sup_{t \in \Theta(M)} |L(t)| \geq z - b_h N_h \right) \\ & \leq \mathbb{P}_f (N_h |L(t_0)| \geq z/2 - b_h N_h/2) \\ & \quad + \mathbb{P}_f \left( N_h \sum_{l=1}^{\infty} \sup_{\substack{u, v \in \Gamma_l \times \Gamma_{l-1} \\ \|u-v\|_1 \leq a^{-l}}} |L(u) - L(v)| \geq z/2 - b_h N_h/2 \right). \end{aligned}$$

We define the function  $\mathcal{W}_t(x, z) = \frac{1}{N_h} q(z + f(x) - f_t(x)) \left( \frac{x-y}{h} \right)^{\vec{p}} \mathbb{I}_{x \in V_h(y)}$ , then the process  $N_h L(\cdot)$  can be written as an empirical process (sum of independent zero-mean random variables)

$$(5.5.6) \quad N_h L(t) = \sum_{i=1}^n \mathcal{W}_t(X_i, \xi_i) - \mathbb{E}_f \mathcal{W}_t(X_i, \xi_i), \quad t \in \Theta(M)$$

At a fixed point  $t_0$ , we can use classical exponential inequalities for empirical process. If the design is deterministic, Hoeffding inequality is sufficient. But the random design requires the Bernstein inequality (see [Boucheron, Bousquet, and Lugosi \[2004\]](#), for example) which takes into account the variance of random variables. By definition of  $\mathcal{W}_t(\cdot, \cdot)$  (in (5.5.6)), we have

$$(5.5.7) \quad \sum_{i=1}^n \mathbb{E}_f \mathcal{W}_t^2(X_i, \xi_i) \leq 1, \quad \|\mathcal{W}_t(\cdot, \cdot)\|_{\infty} \leq 1/N_h,$$

where  $\|\cdot\|_{\infty}$  is the sup-norm.

For the first term of (5.5.5), we use Bernstein inequality. Indeed, we are in the presence of independent random variables bounded by  $1/N_h$  (see definition of  $\tilde{D}_h^{\vec{p}}$  in (5.4.4)) and with finite variance (see 5.5.7), then

$$(5.5.8) \quad \mathbb{P}_f \left( N_h \left| \tilde{D}_h^{\vec{p}}(t_0) - \mathcal{E}_h^{\vec{p}}(t_0) \right| \geq z/2 - b_h N_h/2 \right) \leq 2 \exp \left\{ -\frac{(z - b_h N_h)^2}{8 + \frac{8}{3N_h}(z - b_h N_h)} \right\}.$$

Secondly, in view of (5.5.5), we write the probability with the following form,

$$\begin{aligned} & \mathbb{P}_f \left( N_h \sum_{l=1}^{\infty} \sup_{\substack{u, v \in \Gamma_l \times \Gamma_{l-1} \\ \|u-v\|_1 \leq a^{-l}}} |L(u) - L(v)| \geq z/2 - b_h N_h/2 \right) \\ & \leq \sum_{l=1}^{\infty} \sum_{\substack{u, v \in \Gamma_l \times \Gamma_{l-1} \\ \|u-v\|_1 \leq a^{-l}}} \mathbb{P}_f \left( N_h \frac{\pi^2}{6} l^2 |L(u) - L(v)| \geq z/2 - b_h N_h/2 \right). \end{aligned}$$

Recall that

$$N_h [L(u) - L(v)] = \sum_{i=1}^n \mathcal{W}_u(X_i, \xi_i) - \mathcal{W}_v(X_i, \xi_i) - \mathbb{E}_f [\mathcal{W}_u(X_i, \xi_i) - \mathcal{W}_v(X_i, \xi_i)],$$

then we have a sum of independent centred random variables with finite variance and bounded. The main property of Huber function is that its derivative  $q$  is Lipschitz, then we have the following assertions.

$$(5.5.9) \quad \sum_{i=1}^n \mathbb{E}_f [\mathcal{W}_u(X_i, \xi_i) - \mathcal{W}_v(X_i, \xi_i)]^2 \leq \|u - v\|_1^2, \quad \|\mathcal{W}_u(\cdot, \cdot) - \mathcal{W}_v(\cdot, \cdot)\|_{\infty} \leq \|u - v\|_1 / N_h,$$

Using Bernstein inequality and last inequalities, we obtain

$$\begin{aligned} & \mathbb{P}_f \left( N_h \sum_{l=1}^{\infty} \sup_{\substack{u, v \in \Gamma_l \times \Gamma_{l-1} \\ \|u-v\|_1 \leq a^{-l}}} |L(u) - L(v)| \geq z/2 - b_h N_h/2 \right) \\ & \leq 2 \sum_{l=1}^{\infty} \sum_{\substack{u, v \in \Gamma_l \times \Gamma_{l-1} \\ \|u-v\|_1 \leq a^{-l}}} \exp \left\{ -\frac{36 \|u - v\|_1^{-1}}{\pi^4 l^4} \frac{(z - b_h N_h)^2}{8 \|u - v\|_1 + \frac{8}{3N_h}(z - b_h N_h)} \right\} \\ (5.5.10) \quad & \leq 2 \sum_{l=1}^{\infty} \#(\Gamma_l) \#(\Gamma_{l-1}) \exp \left\{ -\frac{36 a^l}{\pi^4 l^4} \frac{(z - b_h N_h)^2}{8 + \frac{8}{3N_h}(z - b_h N_h)} \right\}, \end{aligned}$$

where  $\#(\Gamma_l)$  is the cardinal of  $\Gamma_l$ . Recall that  $\min_{l>0} \frac{a^l}{l^4} = 1$ ,  $z \geq (1 \vee b_h N_h)$  and  $\Sigma = 2 + 2 \sum_{l=1}^{\infty} \#(\Gamma_l) \#(\Gamma_{l-1}) \exp \left\{ -\frac{18 a^l}{\pi^4 l^4} \right\} < \infty$ , using (5.5.4), (5.5.5), (5.5.8) and the last

inequality, we have

$$\mathbb{P}_f \left( N_h \sup_{t \in \Theta(M)} \left| \tilde{D}_h^{\vec{p}}(t) - \mathcal{E}_h^{\vec{p}}(t) \right| \geq z \right) \leq \Sigma \exp \left\{ -\frac{(z - b_h N_h)^2}{8 + \frac{8}{3N_h} z} \right\}.$$

The lemma is proved. ■

**Proof of Lemma 16.** Remember that the event  $(G_\delta^h)^c$  can be written as  $(G_\delta^h)^c = \{\check{\theta}(h) \notin \mathcal{B}(\theta, \delta)\}$  and  $\check{\theta}(h)$  and  $\theta$  are respectively the unique solutions of equations  $\tilde{D}_h(\check{\theta}(h)) = 0$  and  $D_h(\theta) = 0$ . We can remark the following inclusion

$$\{\check{\theta}(h) \notin \mathcal{B}(\theta, \delta)\} \subseteq \left\{ \sup_{t \in \Theta(M) \setminus \mathcal{B}(\theta, \delta)} \|\tilde{D}_h(t) - D_h(t)\|_2 \geq \varkappa_\delta \right\}.$$

where  $\varkappa_\delta = \inf_{t \in \Theta(M) \setminus \mathcal{B}(\theta, \delta)} \frac{\|D_h(t)\|_2}{2}$ . In view of Lemma 13,  $\varkappa_\delta$  is strictly positive and does not depend on  $n$ . Applying the last inclusion, we obtain

$$\mathbb{P}_f((G_\delta^h)^c) \leq \sum_{\vec{p} \in \mathcal{S}_b} \mathbb{P}_f \left( N_h \sup_{t \in \Theta(M) \setminus \mathcal{B}(\theta, \delta)} |\tilde{D}_h^{\vec{p}}(t) - D_h^{\vec{p}}(t)| > \frac{N_h \varkappa_\delta}{\sqrt{D_b}} \right)$$

The assumptions on  $n, h$  in Lemma 16 allow to show that

$$\frac{N_h \varkappa_\delta}{2\sqrt{D_b}} \geq (1 \vee b_h N_h).$$

Using Lemma 15, we have

$$\mathbb{P}_f((G_\delta^h)^c) \leq D_b \Sigma \exp \left\{ -\frac{nh^d \varkappa_\delta^2}{4D_b (8 + 4\varkappa_\delta/(3\sqrt{D_b}))} \right\}.$$

The lemma is proved. ■

**Proof of Lemma 17.** Note that for any  $k \geq \kappa + 1$  and by definition of  $\hat{k}$  (5.3.3)

$$\{\hat{k} = k\} = \cup_{l \geq k} \left\{ |\check{f}^{(k-1)}(y) - \check{f}^{(l)}(y)| > C S_n(l) \right\}.$$

Note that  $S_n(l)$  is monotonically increasing in  $l$  and, therefore,

$$\begin{aligned} \{\hat{k} = k\} &\subseteq \left\{ |\check{f}^{(k-1)}(y) - f(y)| > 2^{-1} C S_n(k-1) \right\} \\ &\cup \left[ \cup_{l \geq k} \left\{ |\check{f}^{(l)}(y) - f(y)| > 2^{-1} C S_n(l) \right\} \right]. \end{aligned}$$

We come to the following inequality: for any  $k \geq \kappa + 1$

$$\begin{aligned}
 \mathbb{P} \left( \hat{k} = k, G_\delta^{h_k} \right) &\leq \mathbb{P} \left\{ |\check{f}^{(k-1)}(y) - \check{f}(y)| > 2^{-1} C S_n(k-1), G_\delta^{h_k} \right\} \\
 (5.5.11) \quad &+ \sum_{l \geq k} \mathbb{P} \left\{ |\check{f}^{(l)}(y) - f(y)| > 2^{-1} C S_n(l), G_\delta^{h_k} \right\}.
 \end{aligned}$$

Note that the definition of  $S_n(l)$  yields

$$N_{h_l} S_n(l) \geq [1 + \ln(h_{\max}/h_l)]^{1/2}.$$

Thus, applying Proposition 7 with  $\varepsilon = C[1 + \ln(h_{\max}/h_l)]^{1/2}$ , we obtain  $\forall l \geq k-1$

$$\begin{aligned}
 \mathbb{P} \left\{ |\check{f}^{(l)}(y) - f(y)| > 2^{-1} C S_n(l), G_\delta^{h_k} \right\} &\leq D_b \Sigma [h_{\max}/h_l]^{-2rd} \\
 (5.5.12) \quad &= D_b \Sigma 2^{-2rd l}.
 \end{aligned}$$

Here, we have also used that  $k \geq \kappa + 1$ . We obtain from (5.5.11) and (5.5.12) that  $k \geq \kappa + 1$

$$\mathbb{P} \left( \hat{k} = k, G_\delta^{h_k} \right) \leq J_2 2^{-2(k-1)rd},$$

where  $J_2 = D_b \Sigma (1 + (1 - 2^{-2rd})^{-1})$ . ■

# Bibliography

- Akaike H. Information theory and an extension of the maximum likelihood principle. In *Second International Symposium on Information Theory (Tsahkadsor, 1971)*, pages 267–281. Akadémiai Kiadó, Budapest, 1973. cf. page(s): 37
- Anscombe F.J. The transformation of poisson, binomial and negative-binomial data. *Biometrika*, 35:246–254, 1948. cf. page(s): 21, 59, 74
- Arias-Castro E. et Donoho D. L. Does median filtering truly preserve edges better than linear filtering? *Ann. Statist.*, 37(3):1172–1206, 2009. cf. page(s): 32, 82, 154
- Astola J., Egiazarian K., et Katkovnik V. Adaptive window size image de-noising based on intersection of confidence intervals (ici) rule. *J. Math. Imaging Vis.*, 16(3):223–235, 2002. ISSN 0924-9907. cf. page(s): 40
- Astola J., Egiazarian K., Foi A., et Katkovnik V. From local kernel to nonlocal multiple-model image denoising. *Int. J. Comput. Vision*, 86(1):1–32, 2010. cf. page(s): 40, 82, 154
- Aubert G. et Aujol J. A variational approach to removing multiplicative noise. *SIAM J. Appl. Math.*, 68(4):925–946, 2008. cf. page(s): 114
- Autin F. *Point de vue Maxiset en estimation non paramétrique*. PhD thesis, Université Paris 7, 2004. cf. page(s): 17, 37, 39
- Autin F. Maxiset for density estimation on  $\mathbb{R}$ . *Math. Methods Statist.*, 15(2):123–145, 2006. cf. page(s): 34
- Autin F., Picard D., et Rivoirard V. Large variance Gaussian priors in Bayesian nonparametric estimation: a maxiset approach. *Math. Methods Statist.*, 15(4):349–373, 2006. cf. page(s): 34
- Autin Florent. Maxisets for  $\mu$ -thresholding rules. *TEST*, 17(2):332–349, 2008. cf. page(s): 34
- Barron A., Birgé L., et Massart P. Risk bounds for model selection via penalization. *Probab. Theory Related Fields*, 113(3):301–413, 1999. cf. page(s): 34



- Bergh J. et Löfström J. *Interpolation spaces. An introduction*. Springer-Verlag, Berlin, 1976. Grundlehren der Mathematischen Wissenschaften, No. 223. cf. page(s): 17
- Bhattachar S. et Sundareshan M.K. Super-resolution of sar images using bayesian and convex set-theoretic approaches. *SPIE proceeding series*, 4053:189–198, 2000. cf. page(s): 119
- Birgé L. et Massart P. Gaussian model selection. *J. Eur. Math. Soc. (JEMS)*, 3(3):203–268, 2001. cf. page(s): 37
- Birgé L. et Massart P. Minimum contrast estimators on sieves : exponential bounds and rate of convergences. *Bernoulli*, 4(3):329–375, 1998. cf. page(s): 15
- Borovkov A. Statistique mathématique. *Mir, Moscou*, 1987. cf. page(s): 14
- Boucheron S., Bousquet O., et Lugosi G. Concentration inequalities. *Springer*, 2004. cf. page(s): 15, 50, 51, 63, 151, 171
- Bousquet O. A Bennett concentration inequality and its application to suprema of empirical processes. *C. R. Math. Acad. Sci. Paris*, 334(6):495–500, 2002. cf. page(s): 15, 66
- Brown L. et Low M. Asymptotic equivalence of nonparametric regression and white noise. *Ann. Statist.*, 24(6):2384–2398, 1996. cf. page(s): 25, 35, 42
- Brown L.D., Cai T. Tony, et Zhou H.H. Robust nonparametric estimation via wavelet median regression. *Ann. Statist.*, 36(5):2055–2084, 2008. cf. page(s): 15, 19, 61, 154
- Buades A., Coll B., et Morel J-M. A review of image denoising algorithms, with a new one. *Multiscale Model. Simul.*, 4(2):490–530, 2005. cf. page(s): 40
- Bunea F., Tsybakov A., et Wegkamp M. Sparsity oracle inequalities for the Lasso. *Electron. J. Stat.*, 1:169–194 (electronic), 2007. cf. page(s): 36
- Cai T. Adaptive wavelet estimation: a block thresholding and oracle inequality approach. *Ann. Statist.*, 27(3):898–924, 1999. cf. page(s): 19, 61, 154
- Cavalier L. Nonparametric statistical inverse problems. *Inverse Problems*, 24(3):034004, 19, 2008. cf. page(s): 25
- Cavalier L. et Golubev Y. Risk hull method and regularization by projections of ill-posed inverse problems. *Ann. Statist.*, 34(4):1653–1677, 2006. cf. page(s): 25, 38
- Cavalier L. et Tsybakov A. B. Sharp adaptation for inverse problems with random noise. *Probab. Theory Related Fields*, 123(3):323–354, 2002. cf. page(s): 25

- Cavalier L., Golubev G. K., Picard D., et Tsybakov A. B. Oracle inequalities for inverse problems. *Ann. Statist.*, 30(3):843–874, 2002. Dedicated to the memory of Lucien Le Cam. cf. page(s): 25
- Chang X.W. et Guo Y. Huber’s m-estimation in relative gps positioning: computational aspects. *Journal of Geodesy.*, 2005. cf. page(s): 32, 154
- Chichignoud M. Minimax and minimax adaptive estimation in multiplicative regression : locally bayesian approach. *Revision*, 2010a. cf. page(s): 34, 35, 42, 55, 56, 85, 113, 156
- Chichignoud M. Minimax adaptive estimation via locally bayesian approach. *Preprint*, 2010b. cf. page(s): 35, 42, 55, 85
- Chichignoud M. Adaptive huber m-estimation in nonparametric regression. *Preprint*, 2010c. cf. page(s): 35, 42, 153
- DeVore R. et Lorentz G. *Constructive approximation*, volume 303 of *Grundlehren der Mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 1993. cf. page(s): 17
- Donoho D., Johnstone I., Kerkycharian G., et Picard D. Wavelet shrinkage: asymptopia? *J. Roy. Statist. Soc. Ser. B*, 57(2):301–369, 1995. With discussion and a reply by the authors. cf. page(s): 14, 35, 37, 39, 42
- Donoho D. L. et Johnstone I. M. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3):425–455, 1994. cf. page(s): 34, 37
- Efromovich S. Yu. et Pinsker M. S. A self-training algorithm for nonparametric filtering. *Avtomat. i Telemekh.*, (11):58–65, 1984. cf. page(s): 15, 33, 34, 35, 38, 42
- Ghosal Subhashis, Lember Jüri, et Van der Vaart Aad. Nonparametric Bayesian model selection and averaging. *Electron. J. Stat.*, 2:63–89, 2008. cf. page(s): 54, 81
- Goldenshluger A. et Lepski O.V. Universal pointwise selection rule in multivariate function estimation. *Bernoulli*, 14(3):1150–1190, 2008. cf. page(s): 25, 35, 39, 42, 51, 81, 115, 155, 158
- Goldenshluger A. et Lepski O.V. Structural adaptation via lp-norm oracle inequalities. *Probab. Theory and Related Fields*, 143:41–71, 2009a. cf. page(s): 25, 39, 41, 51, 80, 81, 115, 155
- Goldenshluger A. et Lepski O.V. Uniform bounds for norms of sums of independent random functions. 2009b. cf. page(s): 15, 51, 155
- Goldenshluger A. et Nemirovski A. On spatially adaptive estimation of nonparametric regression. *Math. Methods Statist.*, 6(2):135–170, 1997. cf. page(s): 35, 39, 42, 47

- Golubev Y. et Spokoiny V. Exponential bounds for minimum contrast estimators. *Electron. J. Stat.*, 3:712–746, 2009. cf. page(s): 15
- Hall Peter et Jones M. C. Adaptive  $M$ -estimation in nonparametric regression. *Ann. Statist.*, 18(4):1712–1728, 1990. cf. page(s): 19, 30, 61, 154
- Härdle W. et Tsybakov A.B. Robust nonparametric regression with simultaneous scale curve estimation. *Ann. Statist.*, 16(1):120–135, 1988. cf. page(s): 30, 154
- Härdle W. et Tsybakov A.B. Robust locally adaptive nonparametric regression. In *Data analysis and statistical inference*, pages 127–144. Eul, Bergisch Gladbach, 1992. cf. page(s): 19, 30, 61, 154
- Härdle W., Kerkycharian G., Picard D., et Tsybakov A. B. *Wavelets, approximation, and statistical applications*, volume 129 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1998. ISBN 0-387-98453-4. cf. page(s): 14, 17
- Härdle W. et Tsybakov A.B. Local polynomial of the volatility function in nonparametric autoregression. *J. Econometrics*, 1:223–242, 1997. cf. page(s): 114
- Has'minskii R.Z. et Ibragimov I.A. *Statistical Estimation, Asymptotic Theory*. Springer-Verlag, Applications of Mathematics, 1981. cf. page(s): 18, 25, 28, 29, 30, 33, 55, 57, 60, 79, 89, 95, 109, 115, 119, 127, 131
- Huber P. et Ronchetti E. *Robust statistics*. Wiley Series in Probability and Statistics. John Wiley & Sons Inc., Hoboken, NJ, second edition, 2009. cf. page(s): 31, 154
- Huber P.J. Robust estimation of a location parameter. *Ann. Math. Statist.*, 35:73–101, 1964. cf. page(s): 32, 154
- Huber P.J. *Robust statistics*. Wiley, New York., 1981. cf. page(s): 32, 61, 154
- Härdle W., Hart J., Marron J.S., et Tsybakov A.B. Bandwidth choice for average derivative estimation. *journal of the American Statistical Association*, 87:417, 1992. cf. page(s): 151
- Juditsky A. Wavelet estimators: adapting to unknown smoothness. *Math. Methods Statist.*, 6(1):1–25, 1997. cf. page(s): 35, 42
- Juditsky A.B., Lepski O.V., et Tsybakov A.B. Nonparametric estimation of composite functions. *Ann. Statist.*, 37(3):1360–1404, 2009. cf. page(s): 115, 155
- Katkovnik V. *Nonparametric identification and data smoothing*. “Nauka”, Moscow (in Russian), 1985. The method of local approximation. cf. page(s): 26
- Katkovnik V. A new method for varying adaptive bandwidth selection. *IEEE Trans. Image Process.*, 47(9):2567–2571, 1999. cf. page(s): 39, 40, 47, 48

- Katkovnik V. et Spokoiny V. Spatially adaptive estimation via fitted local likelihood techniques. *IEEE Trans. Image Process.*, 56(3):873–886, 2008. cf. page(s): 19, 29, 39, 87, 115, 117, 119, 156
- Kerkyacharian G. et Picard D. Minimax or maxisets? *Bernoulli*, 8(2):219–253, 2002. cf. page(s): 34
- Kerkyacharian G., Lepski O.V., et Picard D. Non linear estimation in anisotropic multi-index denoising. *Probab. Theory and Related Fields*, 121:137–170, 2001. cf. page(s): 25, 39, 51, 80, 81, 155, 158
- Kervrann Ch. et Boulanger J. Optimal spatial adaptation for patch-based image denoising. *IEEE Trans. Image Process.*, 15(10):2866–2878, 2006. cf. page(s): 40
- Klutchnikoff N. *On the adaptive estimation of anisotropic functions*. PhD thesis, Aix-Marseille 1, 2005. cf. page(s): 25, 34, 35, 39, 42, 51, 56, 57, 60, 62, 81, 82, 88, 94, 95, 118, 123, 125, 155, 158, 160
- Korostelëv A. P. et Tsybakov A. B. *Minimax theory of image reconstruction*, volume 82 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1993. ISBN 0-387-94028-6. cf. page(s): 114
- Lecué G. et Mendelson S. Aggregation via empirical risk minimization. *Probab. Theory Related Fields*, 145(3-4):591–613, 2009. cf. page(s): 37
- Ledoux M. On talagrand’s deviation inequalities product measures. *ESAIM: Probability and Statistic*, 1:63–87, 1997. cf. page(s): 15, 51, 66
- Lepski O. V., Mammen E., et Spokoiny V. G. Optimal spatial adaptation to inhomogeneous smoothness: an approach based on kernel estimates with variable bandwidth selectors. *Ann. Statist.*, 25(3):929–947, 1997. ISSN 0090-5364. cf. page(s): 35, 39, 41, 42, 51, 88, 117, 123, 137, 157, 158, 166
- Lepski O.V. On a problem of adaptive estimation in gaussian white noise. *Theory of Probability and its Applications*, 35(3):454–466, 1990. cf. page(s): 33, 34, 35, 37, 38, 41, 42, 56, 88, 89, 117, 118, 119, 157, 158, 160
- Lepski O.V. Asymptotically minimax adaptive estimation i. upper bounds. optimally adaptive estimates. *Theory Probab. Appl.*, 36:682–697, 1991. cf. page(s): 35, 38, 41, 42, 88, 117, 121, 123, 137, 157
- Lepski O.V. Asymptotically minimax adaptive estimation ii. statistical models without optimal adaptation. adaptive estimators. *Theory of Probability and its Applications*, 37: 433–468, 1992a. cf. page(s): 35, 39, 42, 89, 119, 121, 160

- Lepski O.V. On problems of adaptive estimation in white gaussian noise. *In Topic in Nonparametric Estimation*, 12:87–106, 1992b. cf. page(s): 35, 39, 42, 160
- Lepski O.V. et Levit B.Y. Adaptive nonparametric estimation of smooth multivariate functions. *Mathematical methods of statistics*, 1999. cf. page(s): 39
- Lepski O.V. et Spokoiny V.G. Optimal pointwise adaptive methods in nonparametric estimation. *Annals of statistics*, 25(6):2512–2546, 1997. cf. page(s): 35, 42, 88, 118, 157, 160
- Mallows C. Some comments on  $c_p$ . *Technometrics.*, 15:661–675, 1973. cf. page(s): 37
- Marteau C. On the stability of the risk hull method for projection estimators. *J. Statist. Plann. Inference*, 139(6):1821–1835, 2009. cf. page(s): 25, 38
- Massart P. Some applications of concentration inequalities to statistics. *Probability Theory*, volume spécial dédié à Michel Talagrand(2):245–303, 2000. cf. page(s): 15, 51, 66
- Massart P. *Concentration inequalities and model selection*, volume 1896 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007. Lectures from the 33rd Summer School on Probability Theory held in Saint-Flour, July 6–23, 2003, With a foreword by Jean Picard. cf. page(s): 15, 35, 51, 66, 67
- Mathé P. The Lepskiï principle revisited. *Inverse Problems*, 22(3):L11–L15, 2006. cf. page(s): 39
- Meyer Y. *Wavelets and operators*, volume 37 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, Cambridge, 1992. Translated from the 1990 French original by D. H. Salinger. cf. page(s): 17, 37
- Nadaraya E.A. On estimating regression. *Theory of Probability and its Applications*, 9(1): 141–142, 1964. cf. page(s): 14
- Nemirovski A. *Topics in non-parametric statistics*, volume 1738 of *Lecture Notes in Math*. Springer, Berlin, 2000. cf. page(s): 20, 27, 37, 147
- Nussbaum M. Asymptotic equivalence of density estimation and Gaussian white noise. *Ann. Statist.*, 24(6):2399–2430, 1996. cf. page(s): 25
- Parzen E. On the estimation of a probability density function and mode. *Annals of Mathematical Statistics*, 33:1065–1076, 1962. cf. page(s): 14
- Peetre J. *New thoughts on Besov spaces*. Mathematics Department, Duke University, Durham, N.C., 1976. Duke University Mathematics Series, No. 1. cf. page(s): 17

- Petrus P. Robust huber adaptive filter. *IEEE Transactions on Signal Processing.*, 47:1129–1133, 1999. cf. page(s): 32, 154
- Plancade Sandra. Nonparametric estimation of the density of the regression noise. *C. R. Math. Acad. Sci. Paris*, 346(7-8):461–466, 2008. cf. page(s): 36
- Polzehl J. et Spokoiny V. Adaptive weights smoothing with applications to image restoration. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 62(2):335–354, 2000. ISSN 1369-7412. cf. page(s): 39
- Polzehl J. et Spokoiny V. Image denoising: pointwise adaptive approach. *Ann. Statist.*, 31(1):30–57, 2003. cf. page(s): 39
- Polzehl J. et Spokoiny V. Propagation-separation approach for local likelihood estimation. *Probab. Theory Related Fields*, 135(3):335–362, 2006. cf. page(s): 19, 29, 30, 39, 87, 115, 117, 119
- Reiss M., Rozenholc Y., et Cuenod C. Pointwise adaptive estimation for robust and quantile regression. 2009. Source: Arxiv. cf. page(s): 19, 30, 61, 81, 155
- Rigollet Ph. et Tsybakov A. B. Linear and convex aggregation of density estimators. *Math. Methods Statist.*, 16(3):260–280, 2007. cf. page(s): 37
- Rosenblatt M. Remarks on som nonparametric estimates of a density function. *Annals of Mathematical Statistics*, 23:832–837, 1956. cf. page(s): 14
- Rousseeuw P. et Leroy A. *Robust regression and outlier detection*. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. John Wiley & Sons Inc., New York, 1987. cf. page(s): 31, 154
- Rousseeuw P.J. Least median of squares regression. *J. Amer. Statist. Assoc.*, 79(388):871–880, 1984. cf. page(s): 154
- Réfrégier P. *Théorie du bruit et applications en physique*. Hermès Science Publications, 2002. cf. page(s): 119
- Simar L. et Wilson P. Statistical inference in nonparametric frontier models: The state of the art. *Journal of Productivity Analysis*, 13:49–78, 2000. cf. page(s): 24, 114
- Spokoiny V. Estimation of a function with discontinuities via local polynomial fit with an adaptive window choice. *Ann. Statist.*, 26(4):1356–1378, 1998. cf. page(s): 39
- Spokoiny V. et Vial C. Parameter tuning in pointwise adaptation using a propagation approach. *Ann. Statist.*, 37(5B):2783–2807, 2009. cf. page(s): 39, 41

- Spokoiny V. G. Adaptive hypothesis testing using wavelets. *Ann. Statist.*, 24(6):2477–2498, 1996. cf. page(s): 35, 42
- Stein C. M. Estimation of the mean of a multivariate normal distribution. *Ann. Statist.*, 9(6):1135–1151, 1981. cf. page(s): 14, 61, 154
- Stone C. Optimal rates of convergence for nonparametric estimators. *Ann. Statist.*, 8(6):1348–1360, 1980. cf. page(s): 33
- Stone C. Optimal global rates of convergence for nonparametric regression. *Ann. Statist.*, 10(4):1040–1053, 1982. cf. page(s): 15, 33, 38
- Talagrand M. Sharper bounds for Gaussian and empirical processes. *Ann. Probab.*, 22(1):28–76, 1994. cf. page(s): 51, 66
- Talagrand M. Concentration of measure and isoperimetric inequalities in product spaces. *Publications Mathématiques de l’I.H.E.S.*, 81:73–205, 1995. cf. page(s): 15, 66
- Talagrand M. New concentration inequalities in product spaces. *Inventiones Mathematicae*, 126:505–563, 1996a. cf. page(s): 15, 66
- Talagrand M. A new look at independence. *Annals of Probability*, 24:1–34, 1996b. cf. page(s): 15, 66
- Tsybakov A. B. *Introduction to Nonparametric Estimation*. Springer Publishing Company, Incorporated, 2008. cf. page(s): 19, 20, 26, 31, 32, 34, 82, 93, 109, 151, 159
- Tsybakov A. B. Robust estimates of a function. *Problems of Information Transmission*, 18(3):39–52, 1982b. cf. page(s): 30, 61, 154
- Tsybakov A. B. Convergence of nonparametric robust algorithms of reconstruction of functions. *Automation and Remote Control*, (12):66–76, 1983. cf. page(s): 30, 61, 154
- Tsybakov A. B. Robust reconstruction of functions by a local approximation method. *Problems of Information Transmission*, 22(2):69–84, 1986. cf. page(s): 30, 61, 154
- Tsybakov A.B. Nonparametric signal estimation when there is incomplete information on the noise distribution. *Problems of Information Transmission*, 18(2):116–130, 1982a. cf. page(s): 30, 32, 61, 154, 159
- Tsybakov A.B. Pointwise and sup-norm sharp adaptive estimation of function on the sobolev classes. *Annals of statistics*, 26(6):2420–2469, 1998. cf. page(s): 35, 42, 56
- Van de Geer S. *Applications of empirical process theory*, volume 6 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 2000. cf. page(s): 64, 154, 165

Van der Vaart A. W. et Van Zanten J. H. Adaptive Bayesian estimation using a Gaussian random field with inverse gamma bandwidth. *Ann. Statist.*, 37(5B):2655–2675, 2009. cf. page(s): 54, 81

Van der Vaart Aad W. et Wellner Jon A. *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996. With applications to statistics. cf. page(s): 67

Watson G. Smooth regression analysis. *Sankhya, Ser.A*, 26:359–372, 1964. cf. page(s): 14

Woodroffe M. On choosing a delta-sequence. *Ann. Math. Statist.*, 41:1665–1671, 1970. cf. page(s): 19







## Résumé

On se place dans le cadre de l'estimation non paramétrique dans le modèle de régression. Dans un premier temps, on dispose des observations  $Y$  dont la densité  $g$  est connue et dépend d'une fonction de régression  $f(X)$  inconnue. Dans cette thèse, cette fonction est supposée régulière, i.e. appartenant à une boule de Hölder. Le but est d'estimer la fonction  $f$  à un point  $y$  (estimation ponctuelle). Pour cela, nous développons un estimateur local de type *bayésien*, construit à partir de la densité  $g$  des observations. Nous proposons une procédure adaptative s'appuyant sur la méthode de Lepski, qui permet de construire un estimateur adaptatif choisi dans la famille des estimateurs bayésiens locaux indexés par la fenêtre. Sous certaines hypothèses suffisantes sur la densité  $g$ , notre estimateur atteint la vitesse adaptative optimale (en un certain sens). En outre, nous constatons que dans certains modèles, l'estimateur bayésien est plus performant que les estimateurs linéaires.

Ensuite, une autre approche est considérée. Nous nous plaçons dans le modèle de régression additive, où la densité du bruit est inconnue, mais supposée symétrique. Dans ce cadre, nous développons un estimateur dit de *Huber* reposant sur l'idée de la médiane. Cet estimateur permet d'estimer la fonction de régression, quelque soit la densité du bruit additif (par exemple, densité gaussienne ou densité de Cauchy). Avec la méthode de Lepski, nous sélectionnons un estimateur qui atteint la vitesse adaptative classique des estimateurs linéaires sur les espaces de Hölder.

**Mots Clés :** Approche Bayésienne Locale, Critère de Huber, Estimation Robuste, Méthode de Lepski, Adaptation Minimax, Régression Multiplicative, Sélection de la Fenêtre, Vitesses Adaptatives Optimales.

## Abstract

We work in the context of nonparametric estimation in the regression model. Firstly, we consider observations  $Y$  where the density  $g$  is known and depends on a regression function  $f(X)$  unknown. In this thesis, this function is assumed regular, i.e. belonging to a Hölder ball. The goal is to estimate the function  $f$  to a point  $y$  (pointwise estimation). For it, we develop a *local bayesian estimator*, constructed from the density  $g$  of the observations. We propose an adaptive procedure based on the Lepski's method, which allows to construct an adaptive estimator chosen from the family of *local bayesian estimators* indexed by the bandwidth. Under some sufficient assumptions on the density  $g$ , our estimator achieves the adaptive optimal rate (in a particular sense). In addition, we remark that, in some models, the bayesian estimator is more efficient than linear estimators.

Secondly, another approach is considered. We consider the additive regression model, where the density of the noise is unknown and assumed to be symmetric. In this framework, we develop the so-called *Huber estimator* based on the idea of the median ( $\ell_1$  criterion). This estimator allows to estimate the regression function, for any density of the additive noise (for example : Gaussian density or Cauchy density). With the Lepski's method, we select an estimator that achieves the rate of conventional adaptive linear estimators on Hölder spaces without depends on the symmetric density of the noise.

**Keywords :** Bandwidth Selector, Adaptive Optimal Rate, Huber Criterion, Lepski's Method, Local Bayesian Fitting, Minimax Adaptation, Multiplicative Regression, Robust Estimation.

**AMS 2000 subject classification:** 62G08, 62G20.